

The Pragmatics of Confidence in Perceptual and Value-based Choice

Nils Erik Tomas Folke

Darwin College

9/2017

This thesis is submitted for the degree of Doctor of Philosophy.

Declaration

This dissertation is the result of my own work and includes nothing which is the outcome of work done in collaboration except as declared in this Preface and specified in the text.

It is not substantially the same as any that I have submitted, or, is being concurrently submitted for a degree or diploma or other qualification at the University of Cambridge or any other University or similar institution except as declared in the Preface and specified in the text. I further state that no substantial part of my dissertation has already been submitted, or, is being concurrently submitted for any such degree, diploma or other qualification at the University of Cambridge or any other University or similar institution except as declared in the Preface and specified in the text

It does not exceed the prescribed word limit for the relevant Degree Committee.

Preface

The data from Experiments 1 and 2 in Chapter 3 was collected by Julia Ouzia, and some of the findings from this chapter were published in the journal *Cognition* in 2016 in collaboration with Julia Ouzia, Peter Bright, Roberto Filipi and Benedetto De Martino. I coded the experimental tasks, analysed the data and thought about the conceptual implications of the work. I was also the lead author of the article together with Julia Ouzia. Roberto Filipi, Peter Bright and Benedetto De Martino reviewed the manuscript. Julia Ouzia reported the results of the *Cognition* paper in her PhD dissertation at Angela Ruskin University. I have since then substantially added to the analyses submitted in the *Cognition* article, so the version of this work presented in the Chapter 3 is at this stage my own.

Experiment 3, reported in Chapter 4 was collected by Catrine Jacobsen and the results of that experiment was submitted by Catrine Jacobsen as part of a dissertation for an MSc in Cognitive Science at UCL. The joint and extended analyses of Experiment 3 and Experiment 4 were published in *Nature Human Behaviour* in 2016 collaboration with Catrine Jacobsen, Steven Fleming and Benedetto De Martino. I did all of the analyses for the paper and wrote the article in collaboration with Benedetto De Martino. Steven Fleming and Catrine Jacobsen reviewed the manuscript. I have revised and extended the analyses of the chapter since publication, and the text as presented in Chapter 4 is my own work.

Experiment 5 was devised and designed in collaboration with Tor Tarantola. Annika Boldt and David Omar Perez provided input on the experimental design. I wrote the experimental task, Tor and I both collected the data and he has done the computational modelling associated with the study (not reported here). This experiment will primarily be presented in the Dissertation of Tor Tarantola for a PhD at the University of Cambridge. It is only included here to report a null finding of the relationship between GSF (a novel eye tracking measure I first introduced in the *Nature Human Behaviour* paper) and choice accuracy.

The remaining chapters are solely the result of my own work.

Publications

Portions of this thesis appear in the following publication

1. Folke, T., Ouzia, J., Bright, P., De Martino, B., & Filippi, R. (2016). A bilingual disadvantage in metacognitive processing. *Cognition*, 150, 119-132.
2. Folke, T., Jacobsen, C., Fleming, S. M., & De Martino, B. (2016). Explicit representation of confidence informs future value-based decisions. *Nature Human Behaviour*, 1.

Acknowledgments

First I would like to thank Benedetto de Martino for adopting me at a time when I was not particularly confident that I would finish this PhD. Joining the BDM Lab is certainly one of the smartest decisions I have made during my time in Cambridge. While on the subject of the BDM Lab, I would like to thank Tor Tarantola, Paula Kaanders, Annika Boldt and all of the other members of the lab for making me a better thinker and for making the late evenings in the office more enjoyable than they should have been. I am also grateful to Freddy for always being there and never failing to make me laugh. For the taxing task of editing this dissertation, and for putting up with me during the writing process I am deeply indebted to Andrea Brody-Barre. While all errors of content are my own, any errors of grammar should probably be blamed on her. Finally a reluctant thank you to Kai Ruggeri for convincing me to stay in academia a bit longer than I had intended.

This research was supported from a studentship by the ESRC.

Summary

Humans can often report a subjective sense of confidence in a decision before knowing its outcome. Such confidence judgements are positively correlated to accuracy in perceptual and memory tasks, but the strength of this relationship (known as metacognitive accuracy) differs across people and contexts. Computationally, confidence judgements are believed to relate to the strength of evidence favouring each option, but it has been suggested that confidence also captures information from other sources, such as response time. This thesis explores the pragmatics of confidence: what factors influence confidence judgements, how accurate confidence judgements are, and how they might influence future behaviour.

Our knowledge of the antecedents of confidence is extended by this work in two ways, by introducing novel predictors of confidence and by increasing our understanding of well-known ones. I find that bilinguals have worse metacognitive accuracy than monolinguals. This bilingual deficiency in metacognitive accuracy cannot be explained by response time and stimulus strength, suggesting that there is at least one important predictor of confidence that remains unaccounted for. I introduce such a predictor in a new eye tracking correlate of confidence: Gaze-shift-frequency, the number of saccades between options, negatively predicts subsequent confidence in perceptual and value-based decisions. In the value domain, the total value of the options is shown to positively relate to confidence despite being negatively related to accuracy, the first such dissociation to be recorded, as far as I am aware. The dissertation extends our understanding of response time as a predictor of confidence by showing that it influences confidence more for judgements that are made after a choice, relative to those made simultaneously with the choice. This differential influence of response time explains the higher metacognitive accuracy of sequential confidence reports.

I explore the consequences of confidence judgements in the context of value-based choice. Lower levels of confidence are associated with changes of mind when the same options recur in subsequent trials. To test whether these changes of mind are rational, I approximate choice accuracy in the value domain. I propose a novel method based on the transitivity of the full choice set, so that choices that violate the most transitive ordering of the items can be treated as errors. I find that participants who were more metacognitively accurate showed a decrease in transitivity violations over time. These results extend prior work linking confidence judgements to error correction in the perceptual domain.

List of Figures

Figure 2.1. Graphical Representations of Univariate Linear Regressions	37
Figure 2.2. Illustration of Signal Detection Theory	43
Figure 2.3. ROC Curves.....	46
Figure 2.4. Schematic of a Drift Diffusion Process	48
Figure 3.1. The Trial Structure of The Dot Discrimination Task for Experiment 1	58
Figure 3.2. The Trial Structure of The Dot Discrimination Task in Experiment 2.....	60
Figure 3.3. Comparing First Order Performance of Monolingual and Bilingual Participants	67
Figure 3.4. Predicting first order accuracy from dot difference and response time.....	68
Figure 3.5. Parameter Estimates from the DDM Model	69
Figure 3.6. Raw confidence judgements by monolinguals and bilinguals, as a function of the accuracy of the response	71
Figure 3.7. Difference in Metacognitive Efficiency Between Monolinguals and Bilinguals	73
Figure 3.8. Mratios as a function of mean response time and group affiliation	74
Figure 3.9. Predictors of Confidence	75
Figure 3.10. Predictors of Confidence by Accuracy	76
Figure 3.11. Map of the Non-linear Relationship between Response Time and Dot Difference in Predicting Accuracy and Confidence, Experiment 1	78
Figure 3.12. Map of the Non-linear Relationship Between Response Time and Dot Difference in Predicting Accuracy and Confidence, Experiment 2.....	79
Figure 4.1. The Task Structure of Experiment 5	92
Figure 4.2. Relation Between Confidence and Choice	97
Figure 4.3. Dynamics of Information Sampling.....	99
Figure 4.4. Factors that Contribute to Confidence	101
Figure 4.5. Confidence Predicts Change of Mind	103
Figure 4.6. Link between Confidence and Transitivity	106
Figure 5.1. Experimental Procedure	118
Figure 5.2. The Relationship Between Confidence and Stimulus Strength for Simultaneous and Sequential Responses.....	121
Figure 5.3. Sequential Confidence Responses are Associated with Higher Metacognitive Efficiency.....	122
Figure 5.4. Second-order Accuracy as a Function of Total Response Time.....	125
Figure 5.5. Predictors of First-order Accuracy and Confidence.....	127
Figure 5.6. The Relationship between RT and Confidence.....	130
Figure A1. The Response Distributions for Sequentially and Simultaneously Reported Confidence	179
Figure A2. Relationship Between Interjudgment Times and Confidence, First-order Accuracy, Stimulus strength and Second-order Accuracy.....	180

List of Tables

Table 3.1. Bilingual Participants' Language History Information, Experiment 1	62
Table 3.2. Bilingual Participants' Language History Information, Experiment 2	64
Table 3.3. Descriptive Statistics for Control measures, Experiment 1	65
Table 3.4. Descriptive Statistics for Control Measures, Experiment 2	66
Table 5.1. First-order response comparisons.....	120

List of Experiments

Experiment	Description	Chapter
1	Monolingual and Bilingual participants completed a two-alternative forced-choice (2AFC) perceptual discrimination task. The aim of the task was to select which one of two circles contained more dots. The difference in dots was updated online following a 1-up-2-down staircase procedure. Confidence was collected on a visual analogue scale (VAS) after each choice. Response times (RT) and Interjudgment times (IJT) were unconstrained.	3
2	As Experiment 1 but response times were constrained to 1.5 seconds.	3
3	Hungry participants completed a 2AFC value-based task. In each trial they had to choose which one of two snack items they preferred. After each choice participants rated their confidence in their choice on a VAS. The eye movements of the participants were tracked during the choice task. RT and IJT were unconstrained. After the choice task participants completed a Becker-DeGroot-Marschak (BDM) procedure where they rated the maximum amount they would be willing to pay for each item.	4
4	Hungry participants rated 72 common snack items based on how much they would be willing to pay for them in a BDM procedure. The ratings were then used to create triplets of items based on the participant preferences. Participants then completed a choice task where they selected their preferred option from the triplet; their eye movements were tracked and stimuli presentation was gaze contingent. After each choice confidence was tracked on a VAS. RT and IJT were unconstrained.	4
5	Participants completed a two-armed-bandit task with the aim to maximise their reward. The magnitude of the reward was constant between the bandits but reward rate differed. The reward rates were independent, so for each trial either option could win, both options could win or neither could win. After each choice participants received feedback about the outcomes from both bandits. Participant eye movements were tracked both during the choice phase and the feedback phase. RT were unconstrained.	4
6	Participants completed a 2AFC perceptual discrimination task, similar to Experiment 1. Each participant completed two sessions of the task; in one session confidence judgements were reported simultaneously with choice and in the other they were reported after choice. The order of the sessions were counterbalanced between participants. Participant eye movements were tracked during the trials and the presentation of stimuli was gaze contingent. RT and IJT were unconstrained.	5

Table of Contents

Declaration.....	2
Preface	3
Publications.....	4
Acknowledgments.....	5
Summary.....	6
List of Figures.....	7
List of Tables	8
List of Experiments	9
Table of Contents	10
1. Introduction.....	15
1.1. Summary.....	15
1.2. The Pragmatics of Confidence.....	15
1.3. What Constitutes an Accurate Confidence Judgment?	16
1.4. How is Confidence Computed?.....	19
1.5. Sequential Sampling Models of Decision Making and Confidence Judgments	21
1.6. The Consequences of Confidence.....	24
1.7. The Subfields of Metacognition Research.....	27
1.8. The Reliability of Confidence Accuracy	32
1.9. Remaining Questions	33
2. Methods	35
2.1. Summary.....	35
2.2. Generalised Linear Models	35
2.3. Hierarchical GLM.....	39

2.4. Signal Detection Theory	42
2.5. Extending SDT: ROC Curves and Metacognition	44
2.6. Extending SDT: Sequential Sampling.....	48
2.7. Software and Computational Implementation	51
3. Evidence of a Metacognitive Deficit in Bilinguals.....	53
3.1. Summary.....	53
3.2. Introduction.....	53
3.3. Methods.....	58
3.3.1. The Dot Discrimination Task, Experiment 1	58
3.3.2. Dot Discrimination Task, Experiment 2	59
3.3.3. Materials.....	60
3.3.4. Procedure.....	61
3.3.5. Participants, Experiment 1.....	61
3.3.6. Participants, Experiment 2.....	63
3.3.7. Hierarchical Models	65
3.4. Results.....	65
3.4.1. Control Measures, Experiment 1	65
3.4.2. Control Measures, Experiment 2.....	65
3.4.3. First Order Performance, Model Free Analyses.....	66
3.4.4. First Order Performance, DDM.....	69
3.4.5. Second Order Performance	70
3.4.6. Exploring Potential Non-linear Relationships	76
3.4.7. Control Analyses	80
3.5. Discussion	80
4. Explicit Representations of Confidence Inform Future Value-based Decisions.....	86

4.1. Summary.....	86
4.2. Introduction.....	86
4.3. Methods.....	90
4.3.1. Experimental Procedures, Experiment 3.....	90
4.3.2. Experimental Procedures, Experiment 4.....	90
4.3.3. Experimental Procedures, Experiment 5.....	91
4.3.4. Participants, Experiment 3.....	93
4.3.5. Participants, Experiment 4.....	93
4.3.6. Participants, Experiment 5.....	93
4.3.7. Eye Trackers	93
4.3.8. Preparation of the Eye Data, Experiment 3.....	93
4.3.9. Preparation of the Eye Data, Experiment 4.....	94
4.3.10. Preparation of the Eye Data, Experiment 5.....	94
4.3.11. Hierarchical Models	94
4.3.12. Drift Diffusion Models	94
4.4. Results.....	95
4.4.1. Relation Between Confidence and Choice	95
4.4.2. Dynamics of Information Sampling.....	97
4.4.3. Factors that Contribute to Confidence	100
4.4.4. Confidence Predicts Change of Mind	101
4.4.5. Link Between Confidence and Transitivity	104
4.5. Discussion.....	107
5. The Timing of Confidence Judgments Influences Metacognitive Performance.....	111
5.1. Summary.....	111
5.2. Introduction.....	111

5.3. Methods.....	116
5.3.1. Experimental Procedure.....	116
5.3.2. Participants.....	119
5.3.3. Eye Tracking.....	119
5.3.4. Hierarchical Models	119
4.3.5. Drift Diffusion Models.....	119
5.4. Results.....	120
5.4.1. First Order Choices.....	120
5.4.2. Three-way Interaction Between Confidence, Stimulus Strength and Accuracy does not Depend on Timing of Confidence Judgments.....	121
5.4.3. Metacognitive Efficiency is Higher for Sequentially Reported Confidence than Simultaneously Reported Confidence.....	122
5.4.4. Can Differences in Processing Time Account for Differences in Metacognitive Efficiency?.....	123
5.4.5. Sequentially and Simultaneously Reported Confidence Cause Differences in Sensitivity to the Variables Predicting First-order Accuracy and Confidence.....	126
5.4.6. Differential Sensitivity to Response Times Explains Differences in Metacognitive Efficiency	129
5.5. Discussion	130
6. General Discussion.....	136
6.1. Summary.....	136
6.2. Overview of Findings.....	136
6.3. The Consequences of Confidence.....	137
6.4. The Causes of Confidence.....	141
6.5. The Computation of Confidence	144
6.6. Conclusion	148

References	149
Appendices.....	162
Appendix 1: List of stimuli in Experiment 3:.....	162
Appendix 2: List of stimuli in Experiment 4:.....	165
Appendix 3: GSF Does Not Predict Choice, but Interacts With Stimulus Strength in a Model that does not Include DDT	177
Appendix 4: Confidence Distributions for Simultaneous and Sequential Confidence Judgments in Experiment 6	178
Appendix 5: Interjudgment Times Are Not Associated with Confidence, Accuracy or Stimulus Strength in Experiment 6	179
Appendix 6: Confidence is Influenced by Both Positive and Negative Evidence for Both Sequential and Simultaneous Confidence Judgments in Experiment 6	180

1. Introduction

1.1. Summary

This dissertation is about the role confidence plays in human decision making, where confidence is defined as a self-reported estimate of uncertainty in a choice. The introduction provides an overview of the current theories of confidence and empirical work informing those theories. It covers how accuracy is operationalised in the confidence domain and theories on the computational foundation of confidence, with special focus on the relationship between sequential sampling theories of decision making and confidence judgments. Furthermore, the reliability of confidence judgments across testing sessions and modalities will be discussed as well as the consequences of confidence. This overview will conclude by highlighting gaps in our current knowledge, and how the empirical chapters in this dissertation address some of these gaps.

1.2. The Pragmatics of Confidence

Human life is full of uncertainty. We face uncertainty about future events (what will the weather be like tomorrow? Who will win the next election?), about other people's mental states (what does my manager think about my performance?) and even about our own preferences (would I prefer to watch an action film or a comedy?). In order for us to make effective decisions, this uncertainty needs to be captured, and there is a wealth of evidence to suggest that estimates of uncertainty are encoded in the human brain (e.g. Bach & Dolan, 2012; Beck et al., 2008; Kepecs, Uchida, Zariwala, & Mainen, 2008; Kiani & Shadlen, 2009). Not all of these uncertainty representations are consciously accessible. For example, our brains account for sensory and motor noise when planning and performing actions, without us ever being aware that there is any noise to begin with (Knill & Pouget, 2004; Trommershäuser, Maloney, & Landy, 2003).

However, people are able to self-report the uncertainty in their decisions, in the form of subjective accuracy, for a great variety of tasks (Fleming & Lau, 2014). These accuracy judgments are called confidence judgments in the psychological literature and tend to correlate with actual performance (though exceptions exist and will be discussed later). Confidence judgments come in two broad forms: a probability judgment that a categorical decision is correct (e.g. the "left building is taller than the right building" or "I prefer a Mars bar over a Snickers bar") or an uncertainty estimate around a continuous quantity (e.g. "I believe the building is about 50 meters high" or "I believe this gamble pays off 25% of the time"; Navajas et al., 2017; Pouget,

Drugowitsch, & Kepecs, 2016). Categorical confidence judgments are by far the more studied of the two, and will be the focus of the work presented here.

Confidence is closely related to metacognition, which is defined as the ability of a cognitive system to monitor its own performance (Flavell, 1979; Fleming, Dolan, & Frith, 2012). However, some researchers dislike the term metacognition, as they feel it implies the presence of a higher-order system that monitors the basic decision process, an idea that they find implausible (Kepecs & Mainen, 2014; Kiani, Corthell, & Shadlen, 2014). The presence or absence of such a higher-order system is one of the on-going controversies regarding what computations underpin confidence judgments, which will be explored more in depth later in this introduction. However, because of its widespread usage, “metacognitive judgments” will be used interchangeably with “confidence judgments” in this dissertation, with neither implying anything about the generative process.

Confidence judgments have also been used to investigate consciousness, as confidence judgments are a self-reported measure of expected performance and as such must be consciously accessible (Koriat, 2007; Nelson, 1996). While this is an interesting research area it depends on a set of complex questions relating to the nature of consciousness that is outside of the scope of this dissertation. Instead this work will emphasise the pragmatics of confidence judgments, focusing on three questions: What factors feed into confidence judgments, what computations underpin confidence judgments and what are the benefits of being able to make accurate confidence judgment? Before this last question can be explored “accuracy” must be defined in the context of confidence.

1.3. What Constitutes an Accurate Confidence Judgment?

Imagine that you notice something flying above you at high speed. Maybe it is a bird, maybe it is a plane, and, just maybe, it is Superman. You think you spotted the colours blue and red, so after careful deliberation you decide that you saw Superman. After relaying this story to a friend, they are justifiably sceptical and ask you how confident you are that you saw the Man of Steel?

Confidence researchers call the original classification of what you saw a first-order choice (though first order choices in confidence research are typically about Gabor patches and random dot kinematograms rather than super heroes). The confidence judgment, on the other hand, is referred to as a second-order choice because it is an evaluative judgment of the first choice. Similarly, a process that contributes to the first-order choice is known as a first-order process and a process that contributes to the confidence judgment is known as a second-order process.

At first glance it might seem trivial to assess first order-performance in this example. After all, either you did spot Superman or you did not. However, say that you always assume it is Superman whenever something flies over your head. If we only test your super hero spotting ability once, and you happen to be correct, an experimenter cannot determine if your accuracy is due to chance or if you are actually skilled at identifying high-velocity objects. By contrast, if we repeat this exercise many times (perhaps by showing you videos a various flying objects in a lab) we may notice that you always pick the Superman option. Imagine we repeat the procedure 10 times and 6 times the correct option happens to be a certain caped super hero – this would mean that you were correct 60% of the time, despite the fact that your response criterion is completely independent of the evidence presented to you! In technical language, your choices were independent of the stimulus strength of the presented stimuli. Stimulus strength denotes the strength of the presented stimulus, on the dimensions that matter for the decision, so if the stimulus is a line stimulus strength might refer to the length of the line, the orientation of the line or the colour of the line, depending of the nature of the discrimination one attempts to make. In the case of this example, the stimulus strength would capture the “supermanness” of the flying objects including some combination of colour, size and shape. To distinguish between participants’ sensitivity to the stimulus strength of the stimuli and their response biases, it is common to analyse decisions like these with signal detection theory (SDT; Green & Swets, 1996). Traditional SDT deals with binary choices (A and B) and assumes that for any given trial, one sample of internal evidence is drawn from a normal distribution. The mean of the normal distribution is determined by which of the two options is correct for the trial in question. Because these internal evidence distributions have arbitrary scales, standard deviations are often set to 1 for both distributions for mathematical convenience. The difference in means between when option A is correct and when option B is correct is denoted d' (pronounced dee-prime) and captures the sensitivity of the respondent (or equivalently the difficulty of the trial). The choice is determined by where the randomly drawn evidence sample falls relative to a choice criterion, denoted c , which captures the participant’s preference for either option independent of the evidence. Because the response criterion denotes a respondent’s tendency to pick a specific option independently of the evidence it is also called the response bias, or bias for short.

Assessing the accuracy of the confidence judgments becomes slightly more complicated. If the confidence judgment is expressed on a scale of probability correct, from 50% (guessing) to 100% (certain), it is meaningful to talk about the calibration of the confidence judgment. If the mean confidence judgment corresponds to the actual proportion correct, the respondent is said to be well-calibrated (Baranski & Petrusic, 1994). If the level of confidence systematically differs

between correct and error choices so that the level of the confidence is diagnostic of the accuracy of the choice, the participant is said to have good confidence resolution, or confidence sensitivity. Because none of the experiments in this dissertation contains confidence judgments that are reported on a probability scale, I will ignore confidence calibration, so terms like “confidence accuracy” or “metacognitive accuracy” will here always refer to confidence resolution.

Historically, metacognitive accuracy has been measured by correlating the accuracy of a set of choices (coded as 1 for correct and 0 as incorrect) with their confidence judgments (on any scale where higher numbers signify higher confidence; Maniscalco & Lau, 2012). This approach is problematic because these correlations do not account for confidence bias (how likely people are to respond with high confidence independent of their accuracy), so two people whose confidence judgments are equally sensitive to their performance, but whose biases differ can have different confidence-accuracy correlation coefficients (Masson & Rotello, 2009). Because signal detection theory separately estimates sensitivity and bias of first order judgments, researchers tried to expand the same framework to achieve similar estimates for second order judgments (Kunimoto, Miller, & Pashler, 2001).

However, directly transferring the original SDT mathematics to the second order domain has turned out to be problematic. SDT assumes that internal evidence is drawn from two Gaussian distributions with equal variances (one distribution for when Option A is correct and the other when Option B is correct). These assumptions tend to work well for first-order choices, especially in two alternative forced choice tasks (Fleming & Lau, 2014), but Galvin et al. showed if first-order internal evidence distributions are normal with equal variances (as traditional SDT assumes), the second order distributions have unequal variances and are highly non-normal (Galvin, Podd, Drga, & Whitmore, 2003), which means that changes in metacognitive bias influences the sensitivity measure (Evans & Azzopardi, 2007).

Recently Maniscalco and Lau solved this problem by creating a sensitivity measure for confidence that exists in the same space as the SDT measure of first-order accuracy (Maniscalco & Lau, 2012, 2014): meta- d' . Because meta- d' exist in the same space as d' , the SDT assumptions hold and as a result bias and sensitivity are independent (Barrett, Dienes, & Seth, 2013). How meta- d' is computed is discussed at length in the methods chapter of the dissertation, but the key takeaway is that it captures the implied first-order sensitivity from the confidence judgments for a participant with perfect metacognitive insight. Maniscalco and Lau assumed that meta- d' would always be lower than d' , but the reverse pattern has been observed, and would be expected in

situations where the confidence judgment captures more information about the true state of the world than the choice, for example when additional information has been gained between the decision and the confidence judgment (Fleming & Daw, 2017). Because first-order sensitivity and second order sensitivity are captured on the same scale in the meta-d' framework, it is meaningful to compare them directly, for example as a ratio: $\text{meta-d}'/d'$. This ratio, commonly referred to as the Mratio (pronounced em-ratio), captures a participant's *metacognitive efficiency*, or how good their metacognitive performance is relative to their first-order performance. This measure is useful because confidence judgments tend to be more accurate for easy trials than for hard trials (Sanders, Hangya, & Kepecs, 2016). By extension, a person who finds the first-order task easy will be more metacognitively accurate than a peer who finds the task hard, even if their introspective ability is the same. Consequently, when introspective ability is the theoretical quantity of interest, researchers should either keep first-order performance constant between participants or compute Mratios. The relationship between first-order and second-order performance is discussed extensively in Chapter 3.

1.4. How is Confidence Computed?

So far this introduction has defined confidence as a behaviour that estimates the accuracy of a choice. This section will move forward by discussing various theories of how confidence is computed, focusing on important points of contention. The fundamental question regarding the computation of confidence is how it relates to the computation of choice. This broad question can be subdivided into more specific elements, the first being whether there are any confidence specific computations at all or whether confidence judgments are simply a by-product of the first order decision (Fleming & Daw, 2017). A theory that can account for both confidence and choice in a single process is more parsimonious and should therefore be preferred, all other things being equal. Traditional SDT treats confidence as a function of the distance between the decision variable and the criterion (Treisman & Faulkner, 1984). Dynamic extensions of signal detection theory that account for decision time as well as choice (discussed in the next section) can derive confidence from the state of the decision variable at the time of choice (Kepecs et al., 2008; Vickers & Packer, 1982). Preliminary neural evidence supports the link between self-reported confidence and the state of the decision variable at the time of choice (De Martino, Fleming, Garrett, & Dolan, 2013; Gherman & Philastides, 2015; Kiani & Shadlen, 2009; Komura, Nikkuni, Hirashima, Uetake, & Miyamoto, 2013; Wei & Wang, 2015).

However, multiple dissociations between first-order performance and confidence performance have been observed and must be accounted for in order for single-process theories to be

plausible. For example, Graziano and Sigman found that varying the time window between the stimulus presentation and the response influenced first-order accuracy and second-order accuracy differently (Graziano & Sigman, 2009). They found that confidence tended to decrease with artificially increased response times independent of first-order accuracy and trial difficulty. Similarly, Lau and Passingham showed that decreasing the time a visual stimulus was shown prior to a mask decreased confidence to a much greater extent than it decreased first order performance (Lau & Passingham, 2006). Because shorter time-windows prior to a mask make the stimuli less likely to be consciously perceived, this finding supports the potential relationship between confidence and phenomenological consciousness. This link was further investigated by Vlassova and colleagues who found that unconsciously processed visual information influenced choices but not confidence judgments (Vlassova, Donkin, & Pearson, 2014). Additionally, there are a set of studies that show that brain lesions and reversible disruptions of neural activity influence confidence judgments and first-order choices differently (Fleming, Ryu, Golfinos, & Blackmon, 2014; Rounis, Maniscalco, Rothwell, Passingham, & Lau, 2010). See Fleming and Daw, 2017 for a more extensive list of studies that have found dissociations in first-order and second order performance.

Some of the dissociations between first-order and second-order performance can be explained by additional processing time for the second-order judgment (Baranski & Petrusic, 1998; Moran, Teodorescu, & Usher, 2015; Resulaj, Kiani, Wolpert, & Shadlen, 2009) or additional evidence being made available after the choice (Bronfman et al., 2015; Navajas, Bahrami, & Latham, 2016). However, in my view the balance of evidence currently favour models that suggest that the computations underlying confidence judgments are at least partially independent from the computations underlying choices. Single process accounts are hard-pressed to explain the dissociation between confidence and accuracy reported in the Vlassova study (2014), as the trial length was the same for all trials. It is also difficult to explain why disrupting certain neural networks decreases confidence accuracy but leaves first order accuracy intact, without arguing that the two judgments rely on different brain networks and, by extension, different computations.

Because of the recorded discrepancies between first-order and second-order performance it seems probable that confidence judgments have their own dedicated computational structure, but questions still remain about the relationship between confidence computations and first-order computations. For example, are confidence judgments computed after a choice is made or are they are computed in parallel with the choice? Baranski and Petrusic evaluated these accounts by measuring the interjudgment times between choices and confidence judgments in a perceptual

task (Baranski & Petrusic, 1998). They reasoned that if confidence judgments were computed after the choice, the interjudgment times should be independent of the difficulty of the trial (which does influence first order response times) but relate to the level of confidence, because it is easier to assess strong evidence for an option than weak evidence. Conversely, if the confidence computations were parallel, little additional decision time would be necessary, as the confidence judgment had been computed in unison with the first-order decision. Their results suggest that confidence judgments might be computed sequentially when the original choice is under speed-stress but in parallel otherwise. Fleming and Daw (2017) have recently proposed a model where the internal evidence underlying confidence is drawn from a separate but correlated distribution to the internal evidence that determine choice. They show that by varying the relative width of the two internal evidence distributions and the strength of their correlation, they can account for all the results reviewed here. However, they have failed to demonstrate that the model can theoretically be falsified, so it is possible that it is too flexible to have predictive utility.

1.5. Sequential Sampling Models of Decision Making and Confidence Judgments

As mentioned in the previous section, SDT has been extended to account for response times as well as choices. These extensions do not treat a choice as the consequence of a single draw of evidence from a distribution; instead, choices are modelled as the result of multiple sequential draws from a set of evidence distributions. In this framework each option is associated with an evidence distribution with a mean that captures the strength of the evidence favouring that option and the variance of the evidence distribution captures uncertainty around that value. For each unit of time, one piece of evidence is drawn from each of the distributions, and this process repeats until one of the options reaches a decision-threshold at which point that option is chosen. The distance between the starting point and the threshold captures the speed-accuracy trade-off. The further away the threshold is, the more draws will be required to reach it, increasing response time. Additionally, because the evidence distributions are assumed to be symmetric, more draws increase the probability that the noise cancels out, and so improves the signal-to-noise ratio in the decision process (see Gold & Shadlen, 2007 for a more technical discussion). Because this family of models characterises a single decision as the result of multiple sequential evidence samples, they are commonly referred to as sequential sampling models (Forstmann, Ratcliff, & Wagenmakers, 2016).

The most commonly used sequential sampling model is the drift diffusion model (DDM; Ratcliff, 1978; Ratcliff & McKoon, 2008). The easiest way to think of a DDM is to imagine a particle that floats between two boundaries, representing two mutually exclusive options. For each unit of time the movement of the particle is drawn from a normal distribution, where the standard deviation is conventionally set to 0.1 (Vandekerckhove, Tuerlinckx, & Lee, 2011) and the mean, referred to as the drift rate, captures the dominant option so that the mean is positive if the upper option is dominant and negative if the lower option is dominant). The ratio between the mean and the standard deviation captures how noisy the decision process is. The separation between the boundaries captures the speed-accuracy trade-off, and the starting point of the particle captures an a priori preference for one of the options. The model also contains a non-decision time parameter, which estimates the non-decision portion of the response time (e.g. motor preparation and actual movement). Drift diffusion models are popular because they capture both the choices and the response times of two-alternative forced choice tasks well (Forstmann et al., 2016; Ratcliff & McKoon, 2008; Vandekerckhove et al., 2011), because their parameters all have obvious psychological interpretations, and because they have been shown to be mathematically equivalent to the sequential probability ratio test (Gold & Shadlen, 2007), which is the normative solution to sequential sampling problems with two options (Bogacz, Brown, Moehlis, Holmes, & Cohen, 2006). Additionally, a wide variety of computer packages makes them relatively easy to fit (for a systematic comparison of different fitting programs see Ratcliff & Childers, 2015). DDMs are discussed more in the methods chapter.

A wealth of choice related phenomena can be studied by allowing DDM parameters to vary between participants, between conditions, or both. One intriguing example is the work by Krajbich and Rangel, which investigated how visual attention influenced value-based decision making (Krajbich & Rangel, 2011; Krajbich, Armel, & Rangel, 2010; Lim, O'Doherty, & Rangel, 2011). In a series of experiments they found that participants' visual attention influenced their choices, as if the options they looked at got a boost to their drift rate. Originally the researchers assumed that this fixation bias interacted with the value of the option, so that drift rate received a greater boost the greater the value difference between the options. Since then, however, another lab directly compared the original interactive model with an additive model, so that the boost in drift rate from looking at an item is constant, regardless of item value (Cavanagh, Wiecki, Kochar, & Frank, 2014). Whilst this distinction might seem like a technicality, it is practically important because the additive model implies that the visual attention bias is most powerful when we are choosing between two low-value options but that it gets marginalised when we are comparing options of greater subjective value. The interactive model makes no such distinction

because the fixation bias scales with the value of the options. Because of its practical and theoretical importance I will compare the additive and interactive models of the fixation bias in the experiments that involve eye tracking (Chapters 4 and 5).

One weakness of the DDM is that it cannot account for confidence judgments (Pleskac & Busemeyer, 2010). The reason for this is that the evidence accumulation for each option is perfectly anti-correlated with the other: the closer the particle gets to the upper response threshold the further it gets to the lower response threshold. As a consequence, the uncertainty in the decision-process is perfectly captured by the response time, so, according to the model, confidence should be a monotonic function of response time (Ratcliff & Starns, 2009). However, the data suggest that both response times and stimulus strength influence confidence independently (e.g. De Martino et al., 2013; Kiani et al., 2014; Vickers & Packer, 1982). One way to address this problem is simply to treat the interjudgment time as a continuation of the random walk of the decision particle, so that a given confidence judgment is the result of the same evidence accumulation process that guided choice, just with additional processing time (Pleskac & Busemeyer, 2010). This model is an instantiation of the single-process theories reviewed above. As such, they too benefit of the parsimony of accounting for both choices and confidence judgments from a single evidence stream. However, they also inherit the problems of these theories in that they struggle to explain dissociations in first-order and second order performance that aren't captured by response time. Furthermore, because this model assumes that the same drift diffusion process that cause the first order choice then continues to cause the confidence judgment, it assumes an association between confidence accuracy and interjudgment times just as the simple DDM (correctly) predicts a relationship between choice accuracy and response times. However, these assumptions do not hold for interjudgment times which are often independent of confidence outcomes (see Baranski & Petrusic, 1998, and Chapter 5).

An alternative way to capture confidence from a sequential sampling framework is to relax the DDM assumption that the accumulators for the two options are perfectly anti-correlated. The result can be conceptualised as two partially independent accumulators that “compete” with each other to reach the same decision-bound, a family of sequential sampling models known as race models. The distance between the winning accumulator and the losing accumulator at the time of decision captures the balance of evidence favouring the winning option, and is therefore a measure of confidence (Kepecs et al., 2008; Vickers & Packer, 1982; Wei & Wang, 2015). Some early models postulated a 1-1 relationship between the balance of evidence and confidence judgments, but modern models commonly involve an intermediary step where the distance between the accumulators are read out by some other process that make it accessible to self-

report (De Martino et al., 2013; Insabato, Pannunzi, Rolls, & Deco, 2010; Shimamura, 2000). Race models have also been extended by computing confidence from a combination of balance of evidence at time of choice and response time (Kiani et al., 2014; Van Den Berg et al., 2016). Utilising both RT and balance of evidence at time of choice does not only fit the confidence data better than other race-models, this approach is also theoretically appealing as the probability of being correct can be recovered from the model for any given balance of evidence and response time. While race-models capture confidence in an elegant way, they perform worse for capturing response time distributions and accuracy judgments than a conventional DDM (Ratcliff & Smith, 2004). Finally, simultaneous choice and confidence judgments have been modelled by postulating one accumulator for each option (e.g, “high confidence option 1”, “low confidence option 1”, “high confidence option 2”, etc.) but keeping all options anti-correlated so that for any given time unit, one option moves closer to the decision boundary and all other options moves away by the same amount (Ratcliff & Starns, 2009). This model has the benefit of being an extension of traditional two alternative DDMs, and as such fits response times and accuracies well, and it can account for simultaneous choices and confidence judgments. However, it is not obvious how it translates to confidence judgments that are given on a continuous scale.

1.6. The Consequences of Confidence

The previous sections of the introduction have covered research that suggests that humans have systems devoted specifically to representing confidence, either as a readout of the uncertainty in the first order evidence accumulation (De Martino et al., 2013; Insabato et al., 2010; Shimamura, 2000) or as a process carried out in parallel with the first order choice (Baranski & Petrusic, 1998; Fleming & Daw, 2017). Furthermore, a number of published experiments suggest that even non-human animals can track uncertainty in perceptual decisions (Couchman et al., 2010; Lak et al., 2014; Smith et al., 1995) and in memory judgements (Fujita, 2009). A corpus of studies suggest that these results cannot be accounted for by simple associative mechanisms but imply second-order representations of uncertainty (See Smith, Couchman & Beran, 2014 for a review and Kepecs & Mainen, 2014 for a principled computational and methodological argument). It is then natural to wonder why there would be a system specifically devoted to confidence; what benefits do confidence judgments bestow on an agent that justify the computational cost?

A problem that face most animals is whether to forage from their current environment or whether to search for greener pastures elsewhere. These decisions benefit both from a representation of the mean value of alternatives as well as a representation of the variation around this mean (Kolling, Berhens, Mars, Rushworth, 2012). Such foraging decisions are

included in a broader class of problems referred to as exploration vs exploitation problems in contemporary computational neuroscience (Cohen, McClure & Yu, 2007). In environments with uncertain outcomes internal representations of uncertainty can help an agent determine whether they should exploit the option that currently has the best expected value or explore the alternatives to get a more accurate model of the different options and maximise total reward over time (Boldt, Bundell & De Martino, 2017). There is some evidence in that certain non-human animals use uncertainty representations in this way. Washburn et al. (2006) applied an experimental paradigm where Macaques had to make a new two choice discrimination every 6 trials. On each trial they had the option to pick an uncertain response which gave no physical reward but provided information on the correct answer for that trial. The Monkeys correctly picked the uncertain option more often in the first trial of a new discrimination relative to the following trials, gaining information that helped them maximise subsequent rewards. More broadly, explicitly representing model uncertainty helps the agent accurately price the value of additional information.

Another function of representing uncertainty is that it allows agents to avoid risky decisions. A class of experimental designs that rely on optional opt-out illustrate this. Opt-out tasks are typically structured as typical two alternative choice tasks where respondents get a big reward if they are correct or a time-out punishment if they are incorrect. On some trials they also have the option to opt-out and get a small reward that is independent of the decision, the option to opt out may come after the initial decision is made, or as an independent option in the choice phase. To test whether agents use this option judiciously one can compare their performance on trials where they have the option to opt out, but chose not to (free choice trials), to those where they have no choice (forced choice trials). Rhesus monkeys, rats and humans all have shown better performance for the free choice trials than the forced choice trials, when objective difficulty is controlled for, indicating that they can use internal estimates of choice accuracy to improve their performance (Foote & Crystal, 2007; Hampton, 2001; Kiani & Shadlen, 2009; Koriat & Goldsmith, 1996). In metamemory research this is known as the trade-off between quantity and accuracy, as no opt out typically leads to a higher total amount correct at the cost of a lower proportion correct of the trials attempted. Koriat and Goldsmith (1996) has shown that better internal monitoring enable people to maximise their accuracy without sacrificing quantity, as higher metacognitive accuracy enable people to selective use opt out for likely errors, while leaving correct trials intact. In the context of metamemory, respondents can also manage the trade-off between quantity and accuracy by adapting the granularity of their memory report to their confidence in the veracity of their memory (Koriat, Goldsmith & Pansky, 2000). This

relates more broadly to the link between confidence and precision, people may spread their resources to deal with a broader set of contingencies when they perceive the future to be uncertain.

But representing uncertainty may not only allow agents to avoid errors, it might also allow them to correct them: It has been suggested that the same system which is responsible for confidence judgments also triggers error correction in perceptual choice (Yeung & Summerfield, 2012, 2014). Preliminary support for this idea has come from EEG studies, where multivariate analyses showed that the EEG signature that signals errors also predicted graded confidence judgments for correct trials (Boldt & Yeung, 2015). It has also been suggested that confidence judgments help groups of people coordinate more effectively as the uncertainty associated with various observations or predictions can be shared between group members (Shea et al., 2014). Just as confidence might improve group decision making by ensuring that more reliable reports are weighted more heavily, confidence might improve learning in contexts with multiple information sources by ensuring that we learn more from sources that are consistently predictive of salient outcomes (Meyniel & Dehaene, 2017). This may explain why metacognitive accuracy predicts learning speed when general intelligence is controlled for (Veenman, Elshout, & Meijer, 1997). In fact, confidence can act as a learning signal even in the absence of external feedback, helping people to improve their performance in contexts where the outcome is not immediately observable (Guggenmos, Wilbertz, Hebart, & Sterzer, 2016).

Dehaene and Sigman (2012) have suggested that confidence plays a central role in higher-order planning, where multiple steps have to be completed successfully in order for the agent to reach their desired outcome. For example, confidence might indicate when a problem is hard so that more information should be offloaded onto the environment (Gilbert, 2015), or that more deliberation is required (Keramati, Dezfouli, & Piray, 2011). In fact, Yeung and Summerfield (2012) suggested that confidence might regulate the boundary separation in decision making task, helping respondents strike a good balance between speed and accuracy. In contexts where rewards only follow the successful completion of multiple steps, confidence in one's performance in the first steps dictates the effort put into subsequent steps (Van den Berg, Zylberberg, Kiani, Shadlen, & Wolpert, 2016). Van den Berg and colleagues paired two perceptual decisions, and only rewarded participants if both choices were correct. They fitted a drift diffusion model to the second choices and found that participant confidence in the first choice predicted the boundary separation in the second choice, so participants invested more time in getting the second choice right when they were confident that the first choice was correct. This provides preliminary support of Yeung and Summerfield's (2012) idea that

confidence helps improve performance by influencing the boundary separation of a drift process. There is further work to suggest that confidence help inform time and effort allocation in problem-solving. For example, when we encounter a problem our sense of confidence in the first answer that come to mind may determine whether we report it immediately or invest effort and time in trying to find a different answer (Thompson, Turner & Pennycook, 2011; Thompson et al., 2013). Conversely, the evolution of confidence over time may inform when it is time to stop effortful search, and either give our best guess or avoid responding (Ackerman, 2014).

Note that the empirical work on the consequences of confidence has so far focused on the immediate consequences, be it an immediate reversal of a previous decision (Boldt & Yeung, 2015) or a change in decision criteria relative to a choice immediately prior (Van den Berg et al., 2016), so it is an open question to what extent confidence judgments influences decisions over time. Also, the consequences of confidence judgments have to the best of my knowledge solely been explored in relation to perceptual decisions. Chapter 4 will extend our current understanding of the consequences of confidence by studying more long term consequences of confidence judgments in the value domain.

1.7. The Subfields of Metacognition Research

Confidence has been studied extensively in at least three different areas of cognitive science: perception, memory and value-based decision-making. Confidence was first studied in relation to psychophysics (Peirce & Jastrow, 1884), and most of the work referenced here relates to confidence of perceptual judgments. Perceptual discrimination tasks are useful when studying the neural underpinnings of confidence because we have a relatively good understanding of the neural underpinnings of perceptual (primarily visual) first order judgments (Heekeren, Marrett, Bandettini, & Ungerleider, 2004). Additionally, certain psychophysics tasks (e.g. the random-dot-kinetogram) allow for much better control of what evidence is available to an agent at any given time, as they can be set up so that the evidence at each time point can be treated as an independent draw from a normal distribution (Gold & Shadlen, 2007). Together, our relatively advanced understanding of the neural and computational underpinnings of perceptual judgments makes disentangling the relationship between first-order and second-order networks slightly easier than in other domains (e.g. Fleming, Weil, Nagy, Dolan, & Rees, 2010).

Confidence judgments have also been extensively studied in relation to memory (Koriat, 2007; Nelson & Narens, 1990). This metamemory research has extensively dealt with two contemporary points of contention in the perceptual choice literature, and the insights from metamemory are not widely appreciated there (but see Boldt, De Gardelle & Yeung, 2017 and

Shea et al, 2014). The first point of contention is whether confidence judgements depend on a direct readout of first-order evidence strength (see Section 1.4.). Metamemory researchers have discussed this topic since the 1960s (Hart, 1965) and it is now considered largely resolved in that field (Koriat, 1997). The second point of contention is how confidence judgements change based on when they are elicited in relation to the first-order choice. Finally, the combination of insights relating to the substrates of confidence and the temporal dynamics of confidence have in turn provided a more nuanced understanding of the functional properties of confidence, as metamemory researchers appreciate that the functional significance of confidence judgements change over time (Koriat & Levy-Sadot, 2001). Each of these points will be explored more in depth below.

In metamemory research the idea that confidence judgements (referred to as feelings of knowing judgements: FOK) were the result of a direct readout of internal evidence was first put forward by Hart (1965). In his model prospective confidence judgements captured the strength of the memory trace, which in turn predicted the probability of recall at a later time. This model has the benefit that it parsimoniously explains why prospective confidence judgements tend to be diagnostic of recall (Koriat & Levy-Sadot, 2001). The alternative account is known as the inferential account of confidence. The inferential account suggests that confidence is the result of the weighted combination of a set of cues, that are diagnostic of performance during typical circumstances. The inferential view does not imply that the agent make a conscious inference whenever they report their confidence, just that the computation that is driving confidence judgements do not have direct access to the evidence that drive the decision (Koriat, 1997). The inferential view makes a number of specific predictions (1) Cues drive confidence because they are diagnostic of accuracy under typical circumstances. (2) Metacognitive accuracy is the function of the validity of the available cues, and whether these cues are weighted appropriately. (3) Experimenters can manipulate metacognitive accuracy by changing the diagnostic validity of the cues that feed into confidence or how participants weight them.

This last point, that confidence judgements and accuracy can be manipulated independently from one another to create so called “metacognitive illusions” provides a clear empirical test that can arbitrate between the direct access account and the inferential account of confidence. In the context of metamemory research the inferential account has convincingly won this debate. For a review of manipulations that influence first order and second order accuracy independently see the general discussion of Koriat (1997), here I will just mention a few examples. Benjamin and colleagues (1998) first had participants answer 20 trivia questions. Each question had a single word response. After coming up with an answer participants were asked to rate the probability

that they would remember their answer 20 minutes later without any pointers to remind them. They then completed a 20 minute distractor task followed by 10 minutes of free-recall for the 20 responses. Time of deliberation for the trivia questions was negatively associated with confidence in subsequent recall but positively associated with actual recall, so time of deliberation had the opposite effect on confidence and accuracy. Chandler (1994) showed participants a set of target and non-target pictures (learning phase), after 15 minutes of waiting participants were completing a 2 alternative forced choice task between the target and a previously unencountered distractor. Each choice was followed by a confidence judgement. In the two-alternative forced choice task both the targets and the distractors were scenes of lakes. In the experimental condition the non-target in the learning phase was also a lake but it in the control condition it was not. The experimental condition, where targets and non-targets in the were related in the learning phase, were associated with lower accuracy but higher confidence judgements than the control condition. Finally, Begg et al (1989) found that participants predicted that concrete and common words would be easier to recognise than abstract and rare words, however, whereas concrete words were associated with higher recognition rates, more common words were associated with lower recognition rates. Begg and colleagues hypothesised that both concreteness and prevalence predicted ease of processing at the time of the first encounter, but only concreteness aided memorability, whereas more common words were in fact less memorable. In sum, there is plenty of evidence that confidence can be experimentally manipulated independently of choice accuracy, providing evidence that the computation of confidence is at least partially independent from the processes driving recall and recognition.

To understand how the temporal dynamics of confidence has been explored in memory research it might be instructive to explore a common research paradigm in meta-memory. Participants first learn lists of word pairs, with one word being the pointer word and the other being the target. Participants are subsequently presented with the pointer and asked to recall the target from memory, if they fail to recall the target they may be asked to select the target from a set of options. Confidence in accurate recall might be extracted before studying the list, after studying the list, when viewing the pointer, and after an answer has been submitted. Confidence in accurate recognition may be recorded after participants have failed to recall the target but before they are presented with a set of options. This variety in the timing of confidence judgements together with the inferential theory of confidence made metamemory researchers question how the cues that influence confidence might change over time. For example, Hertzog and colleagues (1990) had participants learn a word list and rate their probability of recall (judgment of learning: JOL), before testing their actual recall after a delay, they repeated this procedure 3

times with the same participants. The JOL for the first list were predicted by participants' self-rated memory efficacy and were only weakly predictive of memory performance. However, the JOLs of subsequent lists were driven by the performance on the preceding lists and were more accurate, as participants learned to utilise more accurate confidence cues as the task progressed. Another, more dramatic example of the relationship between metacognitive accuracy and the timing of confidence judgements is that the accuracy of JOL improve dramatically if there is delay between the learning session and when the JOLs are recorded (Rhodes & Tauber, 2011). Presumably the cause of this improvement in metacognitive accuracy is that the conditions for the delayed JOLs are more similar to the conditions during recall, so the cues participants have access to in the delayed condition are more diagnostic of actual performance. The hypothesis that delayed JOLs are associated with higher metacognitive accuracy because they are associated with more valid cues is supported by an elegant experiment by Carrol, Nelson and Kirwan (1997). They trained participants on paired word-lists in two within-participant conditions. In the first condition semantically related pairs were trained until participants demonstrated two correct recalls, in the second condition semantically unrelated words were trained until 8 correct recalls. Recall was then tested two or six weeks later. Participants who gave their JOL immediately after the training judged that they were more likely to remember the related, less practiced word pairs, whereas in reality their recall rate was worse for those pairs. However, participants who gave their JOL one day after the training correctly judged that they were more likely to remember the unrelated, more practiced word pairs. In other words, semantic relatedness was a cue that had a strong influence on confidence immediately after training, but was only weakly associated with long term retention, after a day its relative influence of confidence had decreased, resulting in higher metacognitive accuracy.

Finally, the metamemory research has not only examined how the timing of the confidence judgement in relation to the task may influence metacognitive accuracy, but also how the function of confidence may change over time. Koriat and Levy-Sadot (2001) discuss two cues influencing confidence in recall memory, familiarity and accessibility. Familiarity is the extent to which the item that prompts recall is familiar. For example, Reder (1987) first asked participants to rate the frequency of words, and then had them answer general knowledge questions. Some general knowledge questions contained words from the first task, these were associated with higher confidence judgements, but not higher recall or recognition accuracy. Accessibility is how much information comes to mind when searching for the target. For example, Koriat (1993) found that the number of letters people remembered of a target word predicted confidence judgement, regardless of whether those letters were accurate or not. Koriat and Levy-Sadot

(2001) argued that these two cues influence subjective confidence at different times, because they serve different functional roles. Familiarity, in their account, serves as a gating mechanism that informs the agent whether the desired information is likely to be available in memory, and therefore whether it is worthwhile to engage in effortful search. Once search is initiated the amount of information that comes to mind serves as a cue for the probability of finding an accurate response. This cue is informative because the information that comes to mind during effortful search is mostly accurate under the normal conditions (Koriat & Goldsmith, 1994; 1996). Koriat and Levy-Sadot presented two strands of supporting evidence for their model: First, when they manipulated familiarity and accessibility independently they found an interaction effect so that the effect of accessibility was more pronounced when familiarity was high. This result fit the hypothesis that participants only bothered to engage in effortful search for trials when high familiarity cues indicated that the search might be worthwhile. Second, when the response time of the confidence judgements were constrained, so participants had less time to engage in effortful search, the interaction between familiarity and accessibility diminished. This work suggests a framework for thinking about confidence where the weighting of various confidence cues depends on how predictive they are of accuracy in typical contexts and what functions they fill in relation to the task. I believe this approach could enrich both the theoretical and empirical work on confidence if it was adapted more widely.

To summarise, there are three ideas from metamemory that the broader metacognition research may benefit from: First is the idea that confidence, rather than being a direct function of first-order evidence strength, is inferred from a set of cues. Second, the relative influence of these confidence cues may vary depending on the timing of the confidence judgement. Third, the function of the confidence judgement may also evolve over time. The first idea, that different cues feed into confidence judgements, will be explored throughout this dissertation, the second idea, that the relative weight of these cues may vary with the timing of confidence judgements will be explored in relation to perceptual decision making in chapter 5. The third idea will not be directly explored in this dissertation but provides an interesting avenue for future research.

Finally, confidence has recently been studied in relation to value-based choice (De Martino et al., 2013; Lebreton, Abitbol, Daunizeau, & Pessiglione, 2015). This research field has a disadvantage relative to the other two in that it is harder to evaluate first-order accuracy and, by extension, second order accuracy. For example, if a participant first shows a preference for a Mars bar over a Snickers bar but subsequently changes their mind, it is hard to determine whether the participant made a mistake or whether they genuinely shifted their preference (ignoring the trivial cases where choices can be explained in terms of fast motor errors). However, studies that focus

on value-judgment have nevertheless made important contributions to the study of confidence judgments, as they have shown that the ventro-medial prefrontal cortex, that has been repeatedly implicated in first-order value computations (Basten, Biele, Heekeren, & Fiebach, 2010; Chib, Rangel, Shimojo, & O'Doherty, 2009; De Martino et al., 2013; Hare, Camerer, & Rangel, 2009) is also involved in confidence judgments (Barron, Garvert, & Behrens, 2015; Lebreton et al., 2015). This gives rise to interesting questions regarding the utility of certainty, and more broadly the relationship between confidence and value.

1.8. The Reliability of Confidence Accuracy

Because confidence has been studied in relation to perceptual judgments, memory judgments, and confidence judgments, it is natural to ask how domain-general metacognitive abilities are. Research over the last decade has started to provide an answer. In an early test of the generalizability of metacognitive accuracy, Song and colleagues compared the metacognitive performance of participants across two visual discrimination tasks, and found strong correlations in metacognitive performance but not in first order-performance, and thus suggested that metacognitive abilities might be domain general (Song et al., 2011). McCurdy and colleagues extended this research by comparing confidence performance in a memory task and a perceptual discrimination task. They found that metacognitive performance was correlated across tasks but that the two types of metacognition appeared to be associated with the morphology of different brain regions, and that the correlation in size between these regions accounted for the correlation in metacognitive performance (McCurdy et al., 2013). Further evidence that confidence in memory judgments and confidence in perceptual judgments might be associated with different networks came from a neuropsychiatric study that found patients with anterior prefrontal lesions showed impaired metacognitive performance for perceptual discrimination but intact metacognitive performance in a memory task (Fleming et al., 2014). The picture gets somewhat complicated by a set of studies from Baird which confirmed that separate brain regions were involved in metacognition relating to memory and metacognition relating to perception (Baird, Cieslak, Smallwood, Grafton, & Schooler, 2015; Baird, Smallwood, Gorgolewski, & Margulies, 2013) but showed no correlation in metacognitive performance between the two domains.

This apparent contradiction might be resolved by recent work by Ais and colleagues, who have conducted the most systematic evaluation of the reliability of metacognitive ability to date (Ais, Zylberberg, Barttfeld, & Sigman, 2016). They compared four different tasks (two visual discrimination tasks, one auditory discrimination task, and one short-term memory task) where each task was repeated multiple times by each participant. They found that metacognitive

performance was highly stable across sessions for the same participant completing the same task (confidence biases and the variances were also highly stable across sessions). They also found that metacognitive ability was correlated across modalities. However, they found that metacognitive ability was less stable when the task structure changed. This last point is important, because McCurdy and colleagues maintained very similar task structures for both the memory task and the visual discrimination task (two alternative forced choice tasks between two options presented on different sides of the screen). The Baird experiments, on the other hand had quite different task structures between the domains. The visual discrimination task sequentially showed two sets of gabor patches and the participants had to determine which of the sequences contained a patch that was slightly tilted relative to the other. For the memory task participants had memorised a set of 160 words prior to the test session, during the test session they were shown one word per trial, and had to indicate whether it was in the test set or not. It is therefore possible that the dissociation in performance recorded by Baird and colleagues might be driven by differences in the task structure between the tasks rather than differences in modality.

To summarize, it seems as if metacognitive accuracy correlates across domains in healthy individuals, even though different forms of metacognitive performance is associated with different networks and can be dissociated when brains are damaged or disordered (David, Bedford, Wiffen, & Gillean, 2012; Metcalfe, Van Snellenberg, DeRosse, Balsam, & Malhotra, 2014).

1.9. Remaining Questions

From this brief review of confidence research it is clear that confidence tracks first-order performance across a number of domains and that the relative strength of this relationship appears to be stable within individuals. But what is driving individual variation in metacognitive accuracy? Findings from empirical Chapter 3 will consider a surprising predictor of metacognitive performance: how many languages you speak. This review also highlighted that first order task structure seems to influence how stable metacognitive accuracy is within the same person. But what aspects of the task structure matters? In Chapter 5 I show that the timing of confidence judgments relative to the choice can influence metacognitive accuracy, even when the first-order task is identical across sessions. I explore the cause of this difference by examining how the relationship between confidence and its predictors change as a function of the timing of the confidence judgment. On a trial-by-trial basis, the strength of a confidence judgment is influenced by at least two elements of the first order decision, stimulus strength and confidence,

but what other components of the first order-decision might influence confidence? In Chapter 4 I introduce two novel predictors of confidence, one relating to the perceived value of the presented stimuli and one relating to the visual behaviours of the participants, and I go on to replicate the relationship between this new behavioural index of uncertainty and confidence in Chapter 5. I also show that the amount of time fixating on an option positively predicts the probability of choosing that option, when the stimulus strength of the option is accounted for, in both perceptual and value-based choice (Chapters 4 & 5). Together, these findings suggest that eye behaviours influence confidence and choice in the same way for perceptual and value-based decisions, suggesting that visual attention plays an important domain-general role in how we construct evidence. Finally, I have reviewed the benefits of accurate confidence judgments such as the possibility to quickly correct errors, but it has previously been unclear if confidence is associated with improved decision making over a longer time scale. I provide tentative evidence that this might be the case in Chapter 4. As such, this dissertation provides a number of novel insights regarding the origins, computational underpinnings, and consequences of confidence judgments.

2. Methods

2.1. Summary

This thesis relies on two mathematical frameworks: generalised linear models (GLM) and signal detection theory (SDT). GLMs are a common set of mechanism-agnostic models that are capable of capturing both linear and non-linear relationships. Because most data in this dissertation is hierarchically structured (trials are structured in participants and conditions) this method section will also cover hierarchical extensions of the GLM framework. Hierarchical models enable researchers to examine individual differences, within-participant effects, and differences between groups in a single model, avoiding the problems associated with ignoring any of these sources of variance. Parameters in the SDT framework can also be estimated hierarchically, but because GLM is a more common and well-known framework, hierarchical modelling is introduced in that context. Traditional SDT models derive independent measures of sensitivity and bias in two-alternative-forced-choice (2AFC) tasks. SDT can be extended to capture the sensitivity of confidence judgments in first order choices. Alternatively it can be extended to account for response times as well as choices by modelling how evidence is repeatedly sampled over time, so-called Sequential Sampling Models (SSM). This dissertation relies on a common class of SSM known as Drift Diffusion Models (DDM). DDMs attempt to capture the internal decision process in 2AFC tasks, by simulating a particle moving between two boundaries. DDMs have successfully captured relationships between stimulus strength response times and accuracy in a number of domains and can be considered a mathematical simplification of neurologically plausible decision processes (e.g. Roxin & Ledberg, 2008; see Introduction).

2.2. Generalised Linear Models

GLM refers to a set of models that are based on linear regression. The simplest generalised linear model is a univariate linear regression. It requires two vectors of quantitative data of equal length: Y the outcome of interest and X the predictor (note that formula in this chapter will use upper case letters to denote vectors and lower case letters to denote specific values). A linear regression is a function that transforms X to Y using the following form:

$$Y = \beta_0 + X\beta_1 + E$$

Where β_0 and β_1 are free parameters that are estimated to satisfy some objective function, commonly the maximum likelihood (Friedman, Hastie, & Tibshirani, 2001). $\beta_0 + X\beta_1$ can be thought of as the predicted values of Y, denoted \hat{Y} . β_0 is commonly referred to as the intercept parameter because it specifies the value of \hat{y} when $x = 0$ (the intercept if you plot the model) and β_1 is referred to as the slope parameter because it specifies how much \hat{y} changes for each unit of x (the slope of the curve if you plot the model; see Figure 2.1. a). E is an error vector that captures the difference between \hat{Y} and Y , termed residual error. It is easy to generalise the univariate case to a model with multiple predictors:

$$Y = \beta_0 + X_1\beta_1 + X_2\beta_2 + \cdots + X_n\beta_n + E$$

Where n signifies the total number of predictors in the model.

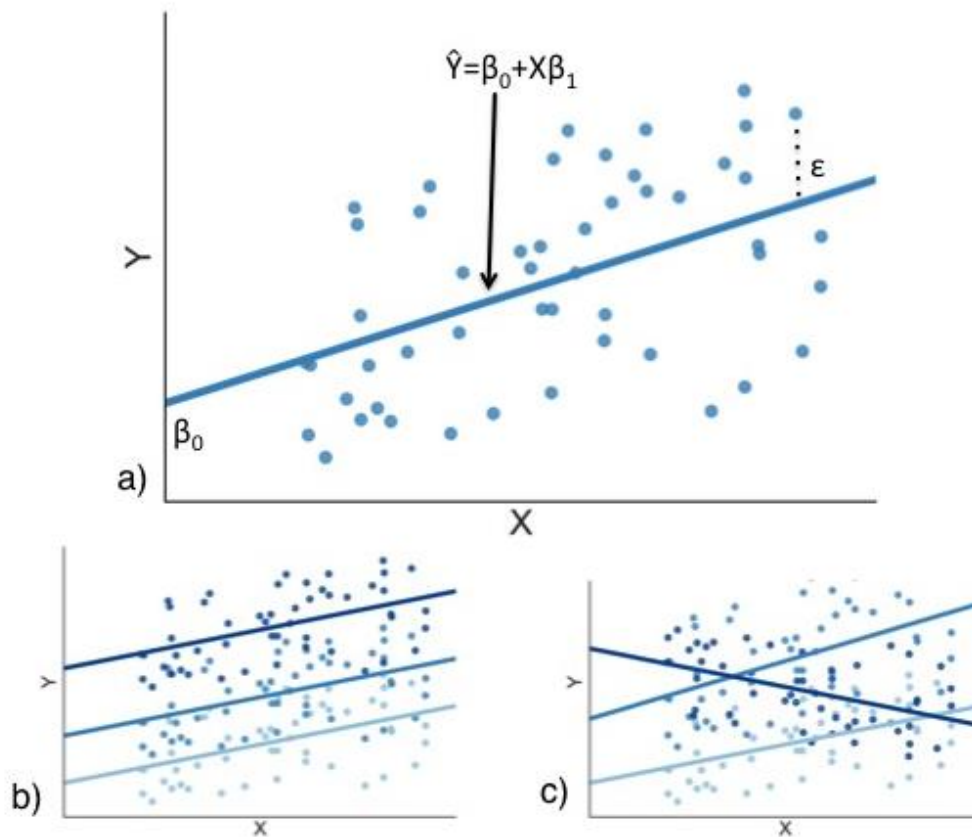


Figure 2.1. Graphical Representations of Univariate Linear Regressions

(a) Shows a basic linear regression, with the mathematical annotation added to the graph. The prediction of a univariate linear model (\hat{y}) is the function of an intercept parameter (β_0) and a predictor (x) multiplied by a slope parameter (β_1). The difference between the prediction \hat{y} and the true value of y for any given x is captured by the error term ϵ . **(b)** and **(c)** Show extensions of linear regression to account for nested data (e.g. when trials are nested in participants). **(b)** Allows for the intercept to vary for each higher ordered structure, **(c)** Allows both the intercepts and slopes to vary. In both **(b)** and **(c)** the varying parameters are constrained in that they are modelled as draws from some distribution.

Linear regression models rely on a number of assumptions to be generalisable and valid (Gelman & Hill, 2006). The first obvious assumption is that the predictor(s) influences the outcome in a linear fashion. However, even if the predictor has a non-linear effect on the outcome, a linear regression model may still be used if the predictor is adequately transformed. For example, if X has a quadratic relationship with Y , $X^2\beta_n$ will capture that relationship even though the model itself is linear. Second, if the model has multiple predictors, the influence of those predictors on the outcome should be additive. If additivity is violated the model can be extended by terms that

capture the interactions between the predictors (e.g. $X_1X_2\beta_n$). Third, the errors in a linear regression should be independent. Fourth, the residual errors have equal variance for all predicted values of Y. Fifth, the residual errors are normally distributed. This last assumption does not influence the estimation of the regression line, but matters when predicting new data from the model (Gelman & Hill, 2006).

The linear regression framework can be extended to binary categorical outcomes (such as picking the left or right option in a 2AFC task) by dummy-coding one outcome as 1 and the other as 0. However, if Y is binary and \hat{Y} is continuous the residual errors cannot be normally distributed, and the variance of the residuals cannot be constant for all \hat{y} , so the assumptions of the regression framework are violated. This problem can be solved by predicting the *probability* of $y_i=1$ (denoted $\Pr(y_i=1)$) rather than predicting Y directly (where i indexes one instance of y in Y). The function that specifies the relationship between \hat{Y} and $\Pr(Y=1)$ is known as a linking function. The most common linking function when predicting a binary outcome from a regression model is the logit function, and a regression using a logit function to predict a binary outcome is called a binomial logistic regression, often abbreviated logistic regression. The logit function is the natural logarithm of the odds. The long form of a univariate logistic regression is:

$$\log \frac{\Pr(y = 1 | x)}{\Pr(y = 0 | x)} = \beta_0 + x_1\beta_1$$

or, equivalently:

$$\Pr(y = 1 | x) = \frac{\exp(\beta_0 + x_1\beta_1)}{1 + \exp(\beta_0 + x_1\beta_1)}$$

Commonly abbreviated as:

$$\Pr(y = 1 | x) = \text{logit}^{-1}(\beta_0 + x_1\beta_1)$$

which is the form I will use for the rest of the thesis.

The logistic regression framework can also be extended to categorical predictions with more than two options by creating a system of logistic regressions that discriminates between each of the options in turn. For example, in a 4-option scenario the first logistic regression would discriminate between whether a participant choose option 1 or any of the remaining three. If they chose one of the remaining options a second logistic regression would discriminate between the second option and the last two. In the univariate case (with options coded as [0, 1, 2, 3]):

$$\Pr(y > 0 \mid x) = \text{logit}^{-1}(\beta_0 + x_1\beta_1)$$

$$\Pr(y > 1 \mid x) = \text{logit}^{-1}(\beta_0 + x_1\beta_1 - \gamma_1)$$

$$\Pr(y > 2 \mid x) = \text{logit}^{-1}(\beta_0 + x_1\beta_1 - \gamma_2)$$

With the constraint that $0 < \gamma_1 < \gamma_2$. The γ variables are called cutpoints and the reason that they are ordered (and all greater than 0) is that the probability of choosing any of the categories must be equal to one. Therefore, if the model estimates that there is a 10% probability of choosing the first option the remaining 3 options must be 90% likely in total, and if the second option is 70% likely to be chosen, the probability of choosing either of the last two options is 20%. This framework can be extended to an arbitrary number of options, though more options require larger samples to get reliable parameter estimates (because the number of free parameters increases with each option).

2.3. Hierarchical GLM

In the real world data is often nested. For example, trials are nested within participants, students are nested within schools, or voters are nested within districts. Models need to account for these patterns to accurately capture real effects. Here I will focus on the case relevant for the empirical chapters, where multiple participants complete many trials inside an experiment.

Say we have a sample of 20 participants completing 500 trials each in a psychophysics experiment and we want to find out whether self-reported confidence increases with response times. Historically, researchers would have either estimated a single regression line for all of the

trials (a pooled model) or analysed each participant individually (an unpooled model). Both of these approaches are suboptimal. The pooled model is suboptimal because it conflates between- and within-participants effects. Returning to our example, suppose that slower participants tend to be more confident. However, within each participant faster trials are associated with higher confidence ratings. By pooling the data we obscure at least one of these effects, depending on which one is dominant, even though there are two very real effects in our data. The problem with the unpooled model is that it is underpowered because it treats people as independent. That is, we model the data as if the first 19 participants provide no information about how participant 20 will behave. Therefore, the unpooled approach needlessly impairs our ability to draw strong conclusions from the data.

There are 3 ways to account for nested data in a GLM framework: allowing the intercepts to vary, allowing the slopes to vary, or allowing both to vary. We will consider a hierarchical model with varying intercepts first.

$$\forall i \in \{1, 2, \dots, n\}$$

$$Y_i = \beta_{0,i} + X_i\beta_1 + E_i$$

The model above estimates one intercept for each participant but only estimates a shared slope parameter for all participants (see Figure 2.1. b). This addresses the risk that effects operating at different levels in the data occlude each other in the model. In terms of the example above, each participant would get their own intercept, capturing that slower respondents were more confident, but there would only be one single slope, capturing that within people faster responses tended to be more confident. However, this approach is still computationally inefficient because it assumes that each intercept is completely independent from the others. This can be addressed by adding a second level to the model, which draws the intercepts at the first level from a distribution of intercepts. The shape of this distribution will influence our results. For convenience, in this example we assume that the intercepts are drawn from a normal distribution:

$$\beta_i \sim N(\mu, \sigma^2)$$

Both levels of the model are fitted simultaneously, so that if all of the participants had similar mean response times, σ^2 would be small, constraining the possible intercepts for outliers. On the other hand, if the mean response time differed a lot between participants, σ^2 would be large, so the hierarchical model would give similar results to the unpooled model. Hierarchical models

have the added benefit that they can be expanded to account for predictors that work on the higher levels in the data structure. In the context of our example, maybe the researchers suspect that participant age influences confidence. This could be tested by adding age (with a slope parameter) into the higher level of the model:

$$\beta_i \sim N(\eta_0 + \eta_1 u_i, \sigma^2)$$

Where u_i indicates the age of the participant, η_1 is a slope parameter and η_0 is an intercept term. The same principle applies if we extend the model to slopes: just as the intercepts above, they can also be allowed to vary by participant, with each slope drawn from a distribution (Figure 2.1. c). Predictors can be added at any level of the model and additional levels can be added (e.g. perhaps we believe that gender influences confidence judgments so we want to nest our participants by gender).

An alternative to maximum likelihood estimation for hierarchical models is Bayesian estimations. In Bayesian estimation the information in the likelihood function is combined with prior information the researcher might have about the data (represented as a distribution). This can be useful when the researcher has a lot of information that is not captured in the data set. However, in cases with large data sets and with diffuse prior distributions, the likelihood tends to dominate the Bayesian computation so that Bayesian and maximum likelihood estimates are very similar. There are 3 important practical advantages with the Bayesian approach: 1) As mentioned above the prior constrains the likelihood, this gives the researcher freedom to implement information from outside of the data in the modelling in a principled way. 2) Bayesian statistics treat all unknown quantities probabilistically so everything from parameter estimates to predictions of new observations are treated as distributions, representing the uncertainty in the model and the data. 3) The software tools implementing Bayesian analysis tends to be more flexible than those implementing pure maximum likelihood approaches so it is easier for researchers to build models that are optimised for their datasets and research questions.

To summarise, the hierarchical GLM framework is both flexible and powerful. Its capacity to evaluate both within-participants and between participants effects simultaneously makes it particularly well-suited to explore questions related to metacognition, where the norm is that a small number of participants complete a large number of trials.

2.4. Signal Detection Theory

Together with GLM, Signal Detection Theory (SDT) is the modelling framework most important for this thesis. SDT translates a noisy observation into a decision. It has been successful because of its relative mathematical simplicity and its explanatory power. It is still relevant today because it can be extended both to capture decisions, confidence judgments and response times. Two extensions are particularly relevant for the work presented in the empirical chapters, the meta-d' framework for measuring confidence accuracy (Maniscalco & Lau, 2012, 2014), and drift diffusion models that capture participant performance in relation to response time and accuracy simultaneously (Ratcliff, 1978). Because both these extensions are highly relevant to the work presented here, I will briefly explain classical SDT as a foundation for the more complex modelling that follows.

SDT is a framework for reaching a decision from a single piece of evidence (traditionally denoted e). In order for e to be meaningful it should carry information about two states we want to discriminate between (H_1 and H_0). For example, say that you try to determine whether it is day (H_1) or night (H_0) only from the outdoors temperature (e). If e is completely determined by the state of H the problem is trivial (if you are in a place where temperature is completely determined by the time of day, you know what time it is as soon as you check the temperature). In interesting real world settings our observations are noisy, so that they are influenced both by the state of H and other factors (e.g. temperature is also influenced by season and weather). To capture this uncertainty e is modelled as a draw from two normal distributions with one mean if H_1 is true and the other mean if H_0 is true (the average temperature is different during the night and during the day, but there is also considerable variation within each of these world states). For mathematical convenience the variance of the distributions are often assumed to be equal, σ^2 . In other words, these normal distributions are the likelihood of $P(e | H_1)$ and $P(e | H_0)$, respectively. The task of the decision maker (and the theory) is then to determine what state of H gave rise to observation e . The difficulty of this decision is determined by the difference in means between H_1 and H_0 scaled by σ ($(\mu_{H1} - \mu_{H0}) / \sigma$), because that difference determines the overlap of the distributions and thus the ambiguity of e (see Figure 2.2.). This difference measure is denoted d' and can be considered the strength of the signal or equivalently the sensitivity of the observer (the day/night discrimination becomes easier in places where temperature does not change much by season and the weather does not vary much and it would be harder if the thermometer is not very precise).

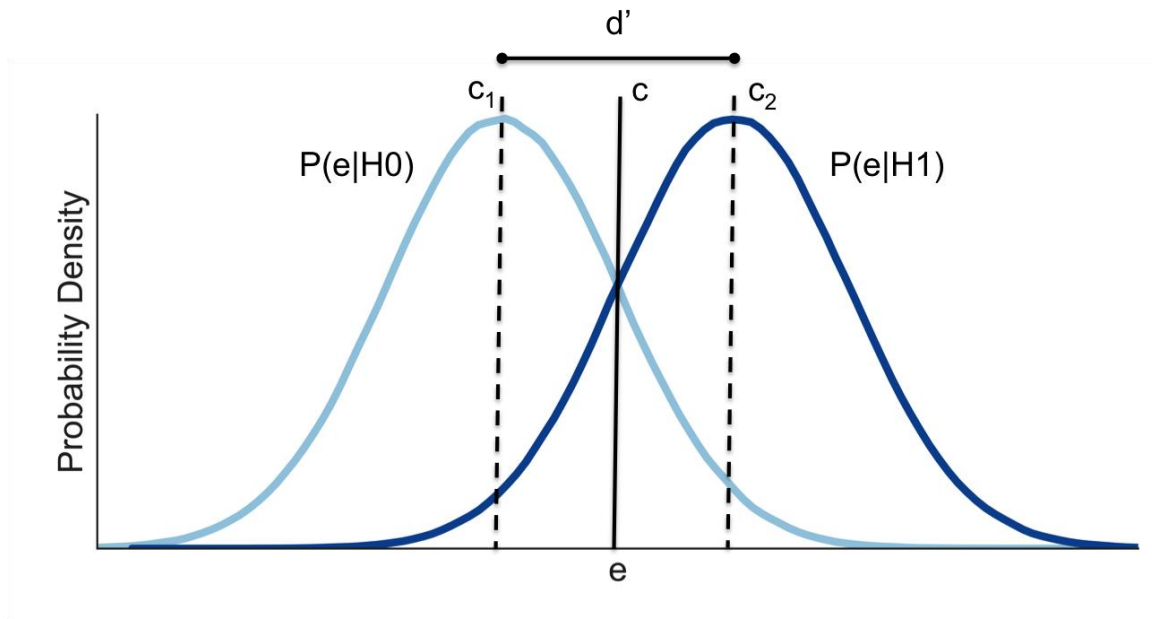


Figure 2.2. Illustration of Signal Detection Theory

According to signal detection theory, the response of a given trial is determined by whether the observation e for that trial exceeds threshold c . e is drawn from one of two Gaussian distributions: the light blue distribution if the target is absent, and the dark blue distribution if the target is present. A participant's discriminative ability is quantified as d' , the difference between the means of the distributions relative to their joint standard deviation. The response threshold c does not have to be in the middle of the distribution but can take any possible value for e . For example, if participants care mostly about hits they may choose C_1 as their criterion, but if they care mostly about avoiding false alarms they might choose C_2 . The figure and text is adapted from Folke, Ouzia, Bright, De Martino, and Filipi (2016).

In order to mathematically discriminate between H_0 and H_1 from e , two things are required: a decision value (DV) and some criterion c that can be compared with the decision value to reach a conclusion. The optimal DV in two choice decisions is the likelihood ratio $l_{1,0}(e) = P(e|H_1)/P(e|H_0)$, or some monotonically related variable. In the temperature example, if we know the probabilities of different temperatures during the day and the night, we can work out the likelihood that it is day or it is night given the temperature we are currently recording. By taking the ratio of these likelihoods we can tell whether it is more likely that it is day or night, and all that is left to do is to pick a decision criterion. When determining what decision criterion to adopt it is useful to consider the four possible outcomes of the decision: either the participant can report H_1 when H_1 is true (a hit), they can report H_1 when H_0 is true (a miss), they report H_0 when H_0 is true (a correct rejection) or they report H_0 when H_1 is true (a false alarm). If the aim is to maximize correct classifications and there are no prior beliefs about the relative

probabilities of H_1 and H_0 c should be 1. If we have different prior beliefs of the respective probability of H_1 and H_0 we can capture that by adopting the response criterion $P(H_1)/P(H_0)$. Finally if the cost of a miss does not equal the cost of a false alarm or the reward of a hit is different from the reward of a correct rejection the optimal c would be given by:

$$\frac{(v_{11} + v_{10})P(H_1)}{(v_{01} + v_{00})P(H_0)}$$

Where v_{ij} is the value of performing action j when hypothesis i is true. Note that when the cost of an error and the reward for a correct response are independent of the state of the world c the equation above simplifies to the prior probability ratio, which lacking any prior information simplifies to 1. Going back to our example, if we only want to be correct as often as possible and we have computed the likelihood ratio $P(\text{temperature} | \text{day})/P(\text{temperature} | \text{night})$ we should behave as if it is day if the ratio is greater than 1 and behave as if it is night, otherwise. However if it is more costly to act as if it is day when it is actually night than it is to act as if it is night when it is actually day, we can modify our response criterion to capture this imbalance.

SDT is of interest to psychologist because the mathematics are not directional, so if a researcher has a record of the hits, misses, correct rejections, and false alarms of a participant they can calculate their sensitivity (d') and response criterion (c). This is useful because it provides two orthogonal and theoretically interesting values that capture most binary choice data well (Stanislaw & Todorov, 1999).

2.5. Extending SDT: ROC Curves and Metacognition

In order to understand how SDT captures actual behaviour and how it can be extended to the metacognitive domain, we need to understand Receiver Operating Characteristic (ROC) curves. ROC curves are graphical non-parametric tools that show how hit rates and false alarm rates are related. Imagine a plane where one axis represents the hit rate, and the other the false alarm rate. We draw a point on this plane that represents the hit rate and false alarm rate of one participant in one particular task (point c in Figure 2.3.). Then we change the incentive structure of the task, so that hits are rewarded more strongly or false alarm are punished less harshly, to motivate the participant to adapt a more liberal response criterion. We capture this new hit rate and false alarm rate with another point on the plane (point c_2 in Figure 2.3.). Once this procedure has been repeated with a variety of different incentive schemes so there is a set of points on the plane, a curve can be drawn to capture the general relationship between hit rates and false alarm rates. The proportion of the plane under this curve is a non-parametric measure of sensitivity (with a

value of 1 for perfect sensitivity and a value of 0.5 for chance level). In signal theoretic terms, this approach incentivises the participant to change their c while keeping d' constant, so we can derive a theoretical ROC curve by mathematically altering c while keeping d' constant and record what that does to hit rates and false alarm rates (see Figure 2.3.). This theoretical curve has been demonstrated to fit the empirical curve well in a variety of different settings (Green & Swets, 1996; Swets, 2014), suggesting that the SDT independence between sensitivity and criterion has empirical support.

There are ways to draw empirical ROC curves without changing the incentive of responses. If participants provide confidence judgments with their choices, other points can be generated in the ROC space by calculating hit rates and false alarms from high-confidence H1 responses, and doing the same with high-confidence H0 responses. Because confidence judgments are second order decisions (i.e. they evaluate the accuracy of the first order decision process) this new curve is called a second-order ROC curve. Interestingly, this second order, confidence-derived ROC curve is different from the ROC curve derived from actually changing the response incentives or drawing a theoretical ROC curve from SDT. Specifically, the area under the curve tends to be smaller when the curve is based on confidence judgments relative to the other methods (see Figure 2.3.).

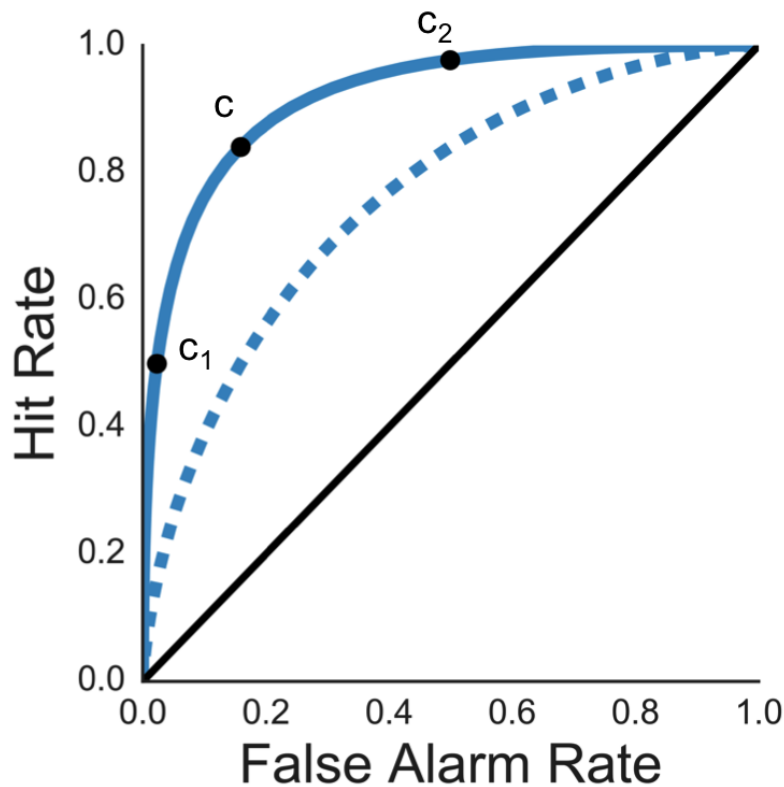


Figure 2.3. ROC Curves

An ROC curve shows the relationship between hit rates and miss rates for a specific sensitivity (d') as the response criterion (c) varies. There are three types of ROC curves that are of interest in this discussion, empirical ROC curves, theoretical ROC curves and confidence ROC curves. Empirical ROC curves are generated by incentivising participants to change their response criterion, recording how that influences the hit rates and false alarm rates (points c_1 , c and c_2 in the figure) and extrapolating a curve from the observations. Theoretical ROC curves are derived by computing d' and c for a single hit rate and false alarm rate and then deriving a curve by keeping d' constant but allowing c to vary. Theoretical and empirical ROC curves tends to overlap (Swets, 2014). Confidence ROC curves can be derived for tasks where confidence judgments on choices are collected. By determining hits and misses based on the confidence judgments, these confidence ROC curves tend to cover a smaller area (the dotted line in the figure) suggesting that some information is lost between the choice and the confidence judgment.

Maniscalco and Lau (2012) invented a way to quantify the difference between the first order-accuracy of confidence judgements and the first-order accuracy of choices. They reversed the

mathematics to determine what level of d' (sensitivity) would have given rise to the second-order ROC curve. This new measure, called meta- d' tells us how well confidence judgments can discriminate between correct and error trials. The calculation of meta- d' is not analytically tractable because it requires extrapolating a curve from the empirical ROC points and then extrapolating a d' from this curve, but there are numeric methods that offer good approximations (see Fleming, 2017; Maniscalco & Lau, 2014). To quantify the difference between first-order sensitivity (the sensitivity based on the choices) and second order sensitivity (sensitivity based on the confidence judgments) Maniscalco and Lau computed the ratio between meta- d' and d' , called the Mratio ($\text{Mratio} = \text{meta-}d'/d'$). An Mratio equal to one means that the confidence judgments are as good at discriminating between H1 and H0 as the first-order responses are, an Mratio greater than one means that they are better, less than one means that they are worse. In most experiments, reported Mratios have been less than 1, suggesting some information loss between the first order and second order judgments (though this is by no means a universal observation, see Fleming & Daw, 2017).

It should be noted that other ways to evaluate the sensitivity of confidence judgments have been proposed. Specifically, high confidence judgments for correct choices could be treated as hits, low confidence judgments for correct choices could be treated as misses, etc. It would then be possible to use these quantities to derive a “second order d' ” and “second order c' ” from the hit rates and false alarm rates of the confidence judgments (Galvin et al., 2003). This approach is hard to implement because traditional SDT assumes that the evidence distributions for both “present” and “absent” trials are normal, and the second order evidence distributions are non-normal when the first order choices are well-captured by these assumptions (Galvin et al., 2003). This results in the undesirable property that metacognitive bias influences the measure of metacognitive sensitivity (Evans & Azzopardi, 2007). Apart from these mathematical challenges, independent modelling work has shown that second order d' is influenced by changes in first order c , whereas meta- d' is stable (Barrett et al., 2013).

Meta- d' can be estimated via maximum likelihood methods or via other forms of analytic approximation (Barrett et al., 2013; Maniscalco & Lau, 2014). Recently Fleming has suggested a new method for hierarchical Bayesian approximation of Meta- d' , this new approach is preferable to the earlier estimation methods for a number of reasons:

1. Point estimates of meta- d' are noisy, Bayesian estimation captures this uncertainty as each meta- d' estimate is represented by a distribution, rather than a point.
2. For research projects where the aim is to compare the metacognitive abilities of two

groups the traditional solution has been to first draw a point estimate for each participant and then compare the samples of point estimates with a t-test. This ignores the uncertainty in the point estimates. A fully Bayesian method would include this uncertainty when estimating the group-level parameters, and conversely use the group level parameters to constrain individual meta-d' estimate, so that noise outliers are pulled towards the group mean.

3. In traditional meta-d' estimation methods padding is used to avoid 0-cells (e.g. if there are 6 levels of confidence ratings, the highest confidence ratings might never have been used for error trials). This padding might bias participant-specific estimates. Bayesian estimation methods do not require padding because the generative multinomial model can handle 0 cell count (Lee, 2008).
4. The general benefits of using Bayesian methods apply to the meta-d' context as well, such as the ability to include information obtained outside the particular experiment (Lee & Wagenmakers, 2014) and evaluating evidence in favour of the null hypothesis (Morey & Rouder, 2011).
5. Simulations have shown that the Bayesian hierarchical meta-d' estimation show better model recovery than the traditional methods, especially in experiments with few trials per participant (Fleming, 2017).

For these reasons, the Bayesian hierarchical approach will be used for the meta-d' analyses presented in the empirical chapters.

2.6. Extending SDT: Sequential Sampling

The flexibility of SDT is both a strength and a weakness. It specifies the relationship between variables but not the absolute values themselves; so as long as the variables scale with each other it is mathematically unconstrained (Gold & Shadlen, 2007). One way to further constrain these models is to account for reaction times as well as decision probabilities. Models that capture both reaction times and choices are called sequential sampling models (SSM), because they are extending SDT to integrate several pieces of evidence before discriminating between H_0 and H_1 . SSM first asks, “Is there enough information to make a successful discrimination?”. If so the appropriate choice is selected; if not another piece of evidence is accumulated and integrated.

The easiest way to visualise a sequential sampling model is to imagine this gradual accumulation of evidence as the movements of a particle (see Figure 2.4.). The starting point of the particle represents our belief before we receive any evidence, then for each unit of evidence the particle moves towards or away from a threshold where the threshold represents a choice. There are

many classes of sequential sampling models (as discussed in the introduction), but I will focus on drift diffusion models (DDM). The distinguishing feature of drift diffusion models is that they measure the relative evidence between two options, so that the evidence accumulation for different options is anti-correlated (the two thresholds are opposite to each other in decision-space, see Figure 2.4.). This can be contrasted with race models, which model evidence accumulation separately for each option and allow for some degree of independence between these accumulators (Forstmann et al., 2016).

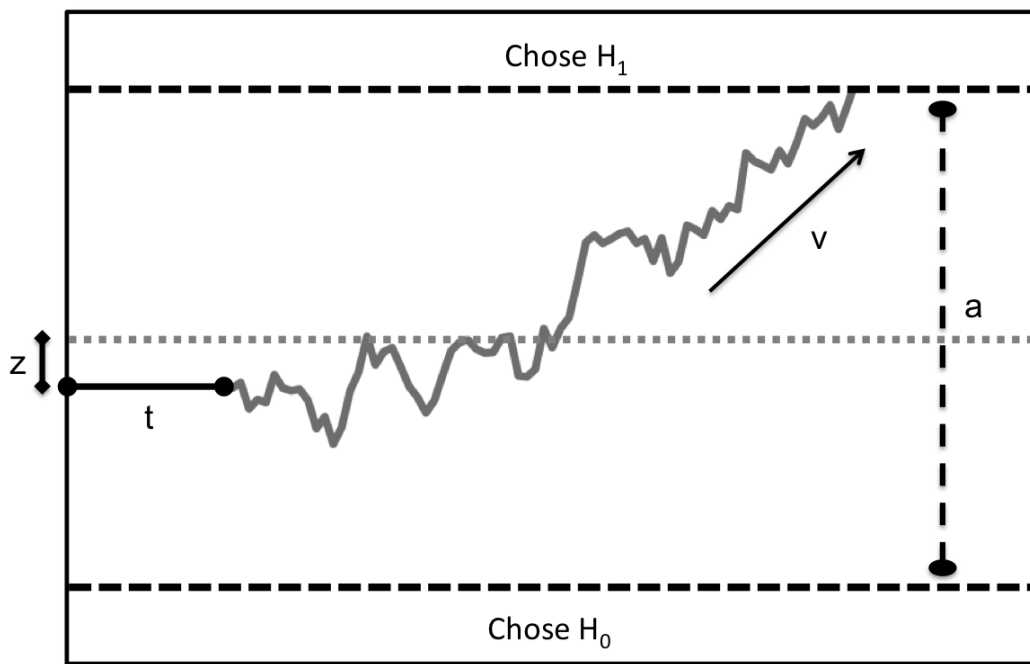


Figure 2.4. Schematic of a Drift Diffusion Process

A drift diffusion process model human decision making as a particle moving between two bounds. 4 key parameters are estimated from the data, bias (z), boundary separation (a) drift rate (v) and non-decision-time (t).

There are 4 primary parameters in a DDM, the drift rate (v) determines the average amount of evidence per unit of time and can be considered a sensitivity parameter (like d' in basic SDT). The boundary separation (a) captures the speed accuracy trade-off; the further away the decision boundaries are the greater the chance that the choice will reflect the drift rate (be accurate) but it will take longer for the particle to hit a boundary, meaning slower responses. As mentioned above, the starting point (z) of the particle captures prior beliefs or preference favouring one option over the other, with the result that choices in the direction of the starting point are more likely, and will happen more quickly when they happen. The boundary separation and starting point together correspond to c in SDT because they jointly capture the response criterion.

Finally, to accurately capture human decision data SDT requires a fourth parameter, non-decision time (t) reflecting processing independent from the decision process itself (such as planning the motor response). In other words response times are a function of the time it takes the evidence to accumulate to a boundary and the non-decision time. Besides the primary free parameters described above there are three additional free parameters: inter-trial variability in drift rates, starting points and non-decision times. These additional parameters are necessary to capture the differences in response time distributions for correct responses and error responses.

It might seem as if DDMs can capture response distributions and response probabilities simply because they have so many free parameters, so any results derived from this approach are psychologically uninformative because the parameters are not sufficiently constrained by the data. However, there are several pieces of evidence to the contrary. First, Ratcliff (2002) used simulations to show that the DDM failed to account for plausible but fake data. Second, the model can be constrained by theoretical considerations (e.g. only allowing task difficulty to influence drift rate, and changing accuracy incentives to influence the speed accuracy trade-off). Third, experiments that vary a specific psychological construct show that most of the resulting variance is being captured by the appropriate parameter (so if difficulty varies between trials, that is captured by the drift rate parameter but not by the threshold parameter; Voss, Rothermund, & Voss, 2004)

The DDM framework has several benefits for behavioural research. Just like SDT, DDM captures both accuracy and bias in an intuitive way, and it also captures the speed accuracy trade-off. This allows researchers to explore questions relating to psychological mechanisms. For example, it is well-known that older people have slower response times for binary decision making tasks than younger people, but is this difference caused by slower evidence integration, slower motor responses, or because they are more careful to avoid errors? DDM analyses show that the differences in response time are best captured by differences in non-decision-times and boundary separation, suggesting that older people have slower motor responses and are more careful, but evidence integration is unaffected by aging (Ratcliff, Thapar, & McKoon, 2010). For a list of research that has benefited from a DDM approach see Forstmann, Ratcliff and Wagenmakers (2016).

The role of experimental manipulations on the DDM parameters can be examined by applying linear regression models to the parameters. For example, if we have two difficulty conditions in a 2AFC experiment and we want to test whether these affect the evidence accumulations we can fit two models. First a null model that ignores trial difficulty when estimating drift-rate and then

a second model that estimates drift rate as an intercept and a coefficient multiplied by a dummy-variable, coded 0 for hard trials and 1 for easy trials. The extent to which the second model fits the data better than the first model would quantify whether the difficulty manipulation has worked, and the size of the coefficient would quantify how much the difficulty manipulation influences evidence accumulation. As with the GLM frameworks we could also draw various intercepts and slopes for different participants, and therefore capture both general trends across people and individual differences between people in the same model.

2.7. Software and Computational Implementation

The hierarchical GLMs in this thesis have been fitted with the lme4 package (version 1.1-7; Bates, Maechler, Bolker, & Walker, 2014) in the R software (version 3.2.3.). P-values and degrees of freedom have been approximated using the Kenward-Rogers approximation, as implemented in the pbkrtest package (version 0.4-2.; Halekoh & Højsgaard, 2014). These approximations are necessary because the complex covariance structures of nested models mean that the null-distributions for the parameters of interest are unknown, and thus need to be approximated. Note that because degrees of freedom are approximated they are continuous variables, rather than integers. Following the suggestion of Gelman and Hill (2006) fixed effect confidence intervals were estimated by multiplying the Wald statistic standard errors by 1.96. Because these confidence intervals are estimates that do not take the covariance between parameters into account (Bolker, 2014) they should not be used to evaluate the significance of a group-level coefficient, but rather serve to give the reader a sense of the precision of the fixed effect coefficients.

Meta-d' and Mratio estimates are obtained through Hierarchical Bayesian estimation, as implemented by Steve Fleming's HMeta-d' code (Fleming, 2017) in MatLab (2016 b). All parameters reported relating to the HMeta-d' code were sampled with three chains, with 11 000 samples in each with the first 1000 samples discarded as burn-in, leaving 30 000 posterior samples for the analyses. All posteriors reported here showed excellent chain mixing (Gelman-Rubin statistics ≥ 1.05), suggesting that the chains had converged and that samples were drawn from the true posterior. All DDM parameters are estimated by Hierarchical Bayesian estimation through Thomas Wiecki and Michael Frank's and HDDM package (version 0.6.0.; http://ski.clps.brown.edu/hddm_docs/) in Python 2.7.12. The HDDM package includes a function to estimate good starting values for the chains prior to starting sampling, dramatically reducing the need for burn-in. All parameters estimated from the HDDM package was sampled with 3 chains, with 2000 samples in each with the first 200 samples discarded as burn-in, leaving

5400 posterior samples for the analyses, the exception being one model that took longer to converge and therefore was sampled with 4000 samples per chain, resulting in 11 400 posterior samples (see Chapter 5). All posteriors reported here showed excellent chain-mixing (Gelman-Rubin statistic > 1.05), suggesting that the chains had converged and that samples were drawn from the true posterior. Some responses with outlying RT's are common in 2AFC tasks, and may provide a serious challenge for likelihood-based DDM estimation (Ratcliff & Tuerlinckx, 2002), therefore the HDDM package offers a mixture model where reaction times are drawn from a mixture of a drift diffusion process and a uniform response distribution that model outliers. In line with the recommendation of the package creators the HDDM models estimated here assumed that approximately 5% of the response times were outliers when estimating the DDM parameters (see the HDDM documentation for more information).

3. Evidence of a Metacognitive Deficit in Bilinguals

3.1. Summary

Some studies have found that bilinguals have an executive function advantage relative to monolinguals, though the exact nature of this advantage, as well as its magnitude, is still uncertain. I tested whether this executive function advantage led to a similar advantage in metacognitive abilities. Two perceptual discrimination experiments showed a metacognitive disadvantage for bilinguals, an effect in the opposite direction than has previously been reported for executive function. I explored how the difference in metacognitive performance relate to how sensitive confidence judgments are to response time and stimulus strength, in the two groups, but fail to find any robust differences. Future work should explore whether these differences are domain general, and if so whether they have any implications for decision making.

3.2. Introduction

Learning a second language brings many benefits: not only does it provide a window into a different culture and enable communication with more people; it might also improve executive function outside the language domain. Early work by Bialystok and colleagues found that bilingual children who used both their languages frequently were better at inhibiting task-irrelevant information than their monolingual peers (Bialystok, 2001). Subsequently they found that this bilingual advantage is driven by the ability to effectively switch between tasks and suppress information relevant to the non-active task (formally known as executive control; Bialystok & Martin, 2004). Other inhibition-related abilities such as response inhibition or delayed gratification do not distinguish bilingual children from their peers (Carlson & Meltzoff, 2008; Martin-Rhee & Bialystok, 2008). The bilingual advantage in executive control has been observed across the life span (Bialystok, Craik, Klein, & Viswanathan, 2004) and in different cultures (Bialystok & Viswanathan, 2009), so it appears to be robust.

The bilingual advantage in executive control has been explained by a theoretical account suggesting that both languages are active by default in the bilingual mind, so the non-active language needs to be constantly suppressed (Green, 1986, 1998). As a result bilinguals practice

executive control almost constantly and get better at it than their monolingual peers. This theory has some empirical support in experimental work on homographs (words that look the same in two languages) which suggests that both languages are active by default for bilingual speakers (Dijkstra, De Bruijn, Schriefers, & Ten Brinke, 2000). Additionally, recent neuroscientific work shows that words from both the target and non-target language are activated during reading, so ignoring words from the non-target language requires active suppression (Van Heuven, Schriefers, Dijkstra, & Hagoort, 2008). There are also a lot of fMRI studies suggesting that the same brain regions are active for both languages, suggesting that there is a single network for language recognition and production and that prefrontal regions modulate which language it is active at any given moment (Abutalebi & Green, 2007). In other words, neuroscience supports the earlier cognitive model which stated that bilinguals have to suppress the non-active language, a designation which constantly shifts depending on context; this constant practice in switching and suppression leads to a bilingual advantage in these domains (Costa, Hernández, & Sebastián-Gallés, 2008). For an up to date theoretical account of how bilingualism is implemented in the brain and how bilingualism might increase executive control see Stocco, Yamasaki, Natalenko, and Prat (2014).

The discovery of bilingual advantages in executive control encouraged research into other differences between bilingual and monolingual cognition. Recent work shows that bilinguals are faster than monolinguals to apply new rules (Stocco & Prat, 2014), they are less prone to egocentric bias when reasoning about other people's beliefs (Rubio-Fernández & Glucksberg, 2012), and they are better able to identify who is speaking and understand what they are saying in the presence of verbal interference (Filippi et al., 2015; Filippi, Leech, Thomas, Green, & Dick, 2012). As such, there might be extensive differences between bilingual and monolingual cognition that are still poorly understood.

However, it is important to note that there have also been studies failing to find differences in executive control between monolinguals and bilinguals (Antón et al., 2014; Gathercole et al., 2014). The earlier positive results have been attributed to researchers failing to match monolinguals and bilinguals on ethnicity and socio-economic status (Morton & Harper, 2007) and to a publication bias favouring positive results (de Bruin, Treccani, & Della Sala, 2015). When Paap and colleagues recently reviewed the state of evidence they concluded that the effect of bilingualism on executive control was either much smaller than previously thought or only manifested during specific and undetermined circumstances (Paap, Johnson, & Sawi, 2015). In sum, it seems like the way in which bilingualism influences executive function is still poorly understood.

Given the interest in the role of bilingualism in cognition in general and executive function in particular it is somewhat surprising that bilinguals and monolinguals have not previously been compared with regards to metacognitive abilities. Metacognition is the ability to monitor and evaluate one's own cognitive processes, or, more informally, to have 'thoughts about thoughts' (Flavell, 1979). Metacognition has been linked to executive function because executive control requires the monitoring of current cognitive processes, and such monitoring is inherently metacognitive (Fernandez-Duque, Baird, & Posner, 2000). Behaviourally, Del Missier and colleagues showed that people who were better at task switching were also better at consistently monitoring their own performance (Del Missier, Mäntylä, & Bruine de Bruin, 2010). On a neurological level the anterior cingulate cortex is involved both with metacognition (Stephen M Fleming & Dolan, 2012) and with conflict monitoring, which is important for executive control (Botvinick, Cohen, & Carter, 2004). With regards to bilingualism, one study found that bilinguals outperformed monolinguals in a task that required carefully monitoring and allocating cognitive resources (Costa, Hernández, Costa-Faidella, & Sebastián-Gallés, 2009). While these results were expressed strictly in terms of executive function, such monitoring processes are inherently metacognitive. There has been some earlier work looking at the role of metacognition in relation to second language learning (e.g. García, Jiménez, & Pearson, 1998; Trofimovich, Isaacs, Kennedy, Saito, & Crowther, 2016), but to the best of my knowledge no prior work has investigated whether monolingual and bilingual people differ with regards to metacognition.

Metacognition can be described as a part of a two-level system, with an object level, or first order process, and a meta level, or second order process (Nelson & Narens, 1994). For example, when choosing which way to drive to work, the object level, first order performance could be operationalised as the travel time of the chosen route relative to the alternatives. Second order performance would be how well our subjective sense of confidence in our choice matches the first order performance; in other words, we would feel more confident if we indeed chose the fastest possible route and less confident if we didn't. The ability to get a subjective sense of one's performance is considered to be a key aspect of confidence judgments (Grimaldi, Lau, & Basso, 2015; C. S. Peirce & Jastrow, 1884). In many cases, subjective confidence judgments are thought to result from an imperfect readout of the uncertainty associated with the first order decisions (Meyniel, Sigman, & Mainen, 2015).

In experimental psychology, metacognitive performance is often assessed by comparing confidence judgments in relation to an objective measure of task performance, such as error rate (e.g., Schwartz & Díaz, 2014; Yeung & Summerfield, 2014). When evaluating metacognitive

performance three terms are of central importance: accuracy, bias, and efficiency (Maniscalco & Lau, 2012, 2014).

Metacognitive accuracy is the extent to which confidence can be used to discriminate between correct trials and error trials (Galvin et al., 2003). For example, if a participant is shown a set of pictures and has to evaluate whether they have seen them before, good metacognitive accuracy would result in their confidence judgments being consistently higher when they are correct compared to when they are wrong. Metacognitive accuracy appears to be domain-general in healthy people, in the sense that people have similar metacognitive accuracy across tasks, even when the tasks depend on different first order abilities (McCurdy et al., 2013; Song et al., 2011; Veenman et al., 1997). However, dissociations in metacognitive abilities have been recorded between memory tasks and visual discrimination tasks (Fleming et al., 2014) in people with localised brain lesions.

In order to gain a complete picture of metacognitive performance one must also account for metacognitive bias, or the tendency to generally report high or low confidence, regardless of the quality of the available information or the accuracy of the first order judgment. For example, people tend to be overconfident in certain memory tasks (i.e., overestimating how often they are correct), whilst still being able to discriminate between correct and incorrect performance (for a review see Hoffrage, 2004). Related to the concept of bias is the concept of calibration. A participant is well calibrated if they pick the correct option 70% of the time, and they are 70% certain they are correct, on average. Calibration is mathematically independent from accuracy because accuracy captures how well confidence can distinguish between correct and incorrect trials whereas calibration captures how well confidence is aligned with performance on average (Baranski & Petrusic, 1994). While questions regarding calibration are interesting, they only really make sense in tasks where confidence is explicitly equated with the probability of being correct, a mapping that does not always feel natural to respondents.

Metacognitive efficiency is a signal theoretic concept that refers to how good a person's metacognitive accuracy is given their first order accuracy. Imagine two people, Susan and John, performing a memory test, followed by a binary confidence judgment. Susan produces fewer errors and therefore has better first order accuracy than John. Nevertheless, both participants report high confidence for 80% of the correctly remembered items and report high confidence for 40% of the items when they were wrong. This means that they both demonstrated the same level of metacognitive accuracy, because their confidence judgments were equally good at discriminating between correct and incorrect trials. However, in a sense John is metacognitively

superior to Susan, because even though his first order decision process is worse, he still shows equally accurate confidence judgments. In the experiments presented below I controlled for first order performance to get a pure measure of metacognitive efficiency in two ways. First, I used an adaptive staircase to ensure a similar first-order accuracy for the experimental task across all participants. Second, I controlled for differences in first order performance mathematically.

Historically, metacognitive accuracy was computed by correlating confidence with first order performance within each participant (Kornell, Son, & Terrace, 2007; Nelson, 1984). However, this approach has been criticised for its inability to distinguish metacognitive accuracy from metacognitive bias (Masson & Rotello, 2009). Maniscalco and Lau (2012, 2014) recently addressed this problem by applying signal detection theory (SDT) to metacognition, thus providing separate measures for bias and sensitivity. Their measure for sensitivity, meta- d' , captures the extent to which confidence judgments can discriminate between correct and error trials. Metacognitive efficiency is captured by the $Mratio$ (Fleming & Lau, 2014), which is computed by dividing meta- d' with d' (a measure of first order sensitivity). This approach to measuring metacognitive performance has been demonstrated to outperform alternatives and to give robust measures of metacognitive accuracy and metacognitive efficiency (Barrett et al., 2013). For a more in depth discussion of these constructs see the Chapter 2.

The aim of this study was to test if the executive control advantage reported in bilinguals translated to a metacognitive advantage. I compared metacognitive efficiency (captured by the $Mratio$) between bilinguals and monolinguals in a perceptual two-alternative-forced choice (2AFC) task. Because metacognitive performance tends to be associated with task performance (Galvin et al., 2003; Maniscalco & Lau, 2012), I used a perceptual task that allowed me to adjust task difficulty online for each participant, titrating performance at around 71% for all participants. This standardisation ensured severely restricted variation in task performance across participants, implying that any variation in metacognitive performance could not be accounted for by differences in task performance. After completing an initial experiment without any response time constraints, I found that the monolingual group responded significantly slower than the bilingual group. Because previous work suggests that metacognitive judgments are partially computed while the first-order decision is made (Baranski & Petrusic, 1998, 2001), this might confound any results relating to metacognition. I addressed this confound in two ways. First, I ran a drift diffusion model on the choice data to test if the differences in response times were best captured by differences in the rate of evidence accumulation, differences in response thresholds, or differences in non-decision-time (for more on drift diffusion models see the methods section). The aim of this model was to see whether either language group accumulated

first-order evidence more efficiently, as such a difference could confound the second-order comparison. Second, I conducted an additional study with a similar set-up, except that participants had to respond within the first 1.5 seconds. In this second experiment response times were comparable across language groups. The results of both experiments are presented together in order to make it easy to compare the results and get a sense of the total findings.

3.3. Methods

3.3.1. *The Dot Discrimination Task, Experiment 1*

Participants completed a two-alternative-forced-choice task programmed in PsychoPy v. 1.82 (Peirce, 2008) presented on a 24-inch widescreen monitor using a standard keyboard. A MATLAB version of a similar task has previously been used in Fleming et al. (2014). On each trial participants saw two white circles on a black background, and indicated whether the left or the right circle contained more dots by pressing the appropriate arrow key on a standard computer keyboard. For every trial, one circle was randomly assigned to have 50 dots; the other circle contained a variable number of dots that was either larger or smaller than 50. The difference in number of dots between the two circles was modified throughout the experiment by a staircase procedure, so that whenever participants correctly responded to two successive trials the task increased in difficulty (one less dot difference between the options) and for every failed trial the task became easier (one more dot difference between the options). The purpose of the staircase was to normalise first order accuracy at 71% across the sample. After each trial participants were asked to indicate their confidence on a sliding scale. Response times were unconstrained for both first order and second order judgments. For a graphical representation of the trial structure see Figure 3.1.



Figure 3.1. The Trial Structure of the Dot Discrimination Task for Experiment 1

62 participants (31 monolinguals and 31 bilinguals) were asked to rate which one of two circles contained more dots and then rated their confidence on a scale going from less to more. Both

the choices and the confidence judgements had unconstrained response-times. One circle always contained 50 dots whereas the other circle contained 50 dots +/- a dot difference. The dot difference was calculated by a 1-up-2-down staircase procedure to fix accuracy at around 70% across participants.

Participants completed 8 blocks with 25 trials in each, making up a total of 200 trials. Prior to beginning the main task, participants completed three practice phases. In the first phase they were shown pairs of circles with the number of dots indicated in writing below the circles. In the second phase participants started making perceptual choices without conducting any confidence judgments. These trials started with a 20-item dot difference, which first changed in increments of four, then in successively smaller increments down to one; this was performed to calibrate the difficulty to each participant. The second phase terminated after 8 reversals (i.e. when participants had switched between picking the correct and the incorrect option 8 times). Participants received feedback on their choices in the second calibration phase. The final phase consisted of 10 trials that simulated the main experimental trials in every way, i.e., without performance feedback, and they were asked to indicate their confidence in their choice after each trial. The edges of the confidence scale was labelled “less” and “more” and participants were instructed to report their relative confidence, i.e. reporting higher confidence on trials they were more certain of being correct and lower confidence when they were less certain of being correct in relation to other trials in the same experiment. They were encouraged to use the full range of the scale. No feedback was given for their use of the confidence scale. All practice trials were excluded from all analyses.

3.3.2. Dot Discrimination Task, Experiment 2

The dot discrimination task was identical to the task in Experiment 1, with the exception that participants now had to respond within 1.5 seconds after first seeing the dots. I also introduced slightly longer inter-trial intervals which featured a fixation cross in the centre of the screen (see Figure 3.2.). If participants took longer than 1.5 seconds to respond, the trial was terminated and the words “Too Slow” were presented for one second.

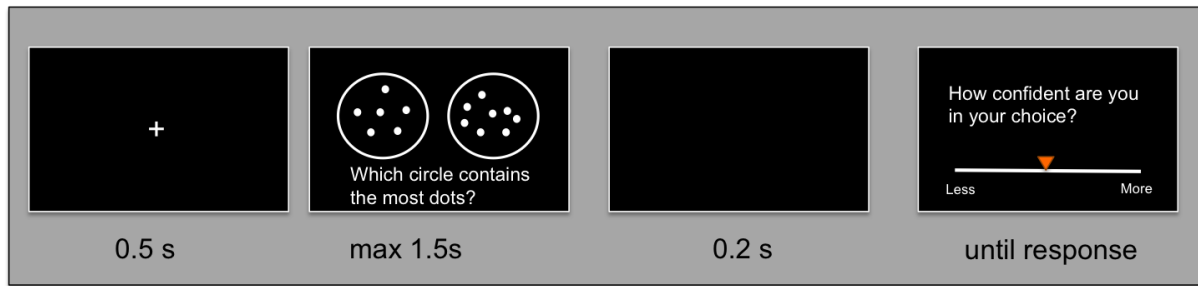


Figure 3.2. The Trial Structure of the Dot Discrimination Task in Experiment 2

60 participants (32 monolinguals and 28 bilinguals) were asked to rate which one of two circles contained more dots and then rated their confidence on a scale going from less to more. Choices were constrained to be faster than 1.5 ms. If participants took longer than 1.5. ms to respond the words “Too Slow” were presented for 1 second before the next trial started. Confidence judgements had unconstrained response-times. One circle always contained 50 dots whereas the other circle contained 50 dots +/- a dot difference. The dot difference was calculated by a 1-up-2-down staircase procedure to fix accuracy at around 70% across participants.

3.3.3. Materials

We administered standardised measures of working memory and non-verbal reasoning to all participants in order to ensure that the groups were comparable with regard to general cognitive function. The materials used were identical for both experiments.

3.3.3.1. Working Memory Test

Working memory was assessed using the digit span task of the Wechsler Adult Intelligence Scale IV (Wechsler, 2008). In this task, participants are asked to repeat a set of single digits (between two and nine) after the experimenter. During the first round (eight sets of two trials) they are asked to repeat the numbers in the same order; in the second round (seven sets of two trials) they have to repeat the numbers in reverse. Each round is terminated once a participant has failed to correctly repeat both trials of one set, and a total score is calculated with a maximum of thirty points.

3.3.3.2. Non-verbal reasoning

Non-verbal reasoning was measured using the Raven’s Advanced Progressive Matrices (Raven & Court, 1986). In this task participants were presented with twelve trials. In each trial they were

shown an incomplete matrix of black and white abstract figures. Participants were asked to identify the missing piece from a selection of eight alternatives and complete all 12 trials in no more than 10 minutes. None of the participants reached this time limit.

3.3.3.4. English language proficiency

In addition to the language history questionnaire, we also measured English language proficiency in bilinguals using the picture naming scale of the Bilingual Verbal Ability Tests (BVAT; Muñoz-Sandoval, Cummins, Alvarado, & Ruef, 1998).

3.3.4. Procedure

The procedures were equivalent for both experiments, except for which version of the dot discrimination task was used. All participants were tested in one hour-long session at Anglia Ruskin University, Department of Psychology in the same room using the same equipment. After informed consent was given they completed a short demographics questionnaire and the bilingual participants were also asked to complete an adapted version of the language history questionnaire by Li, Sepanski and Zhao (2006). We then administered the digit span task, Raven's Matrices as well as the dot discrimination task, with task ordering counter-balanced across participants. The practice blocks of the dot discrimination task were presented with extensive instructions and participants were encouraged to ask any questions prior to task commencement.

3.3.5. Participants, Experiment 1

We tested sixty-two healthy young adults, thirty-one English monolinguals ($M_{\text{age}} = 22.3$, $SD = 3.7$; 12 males), and thirty-one bilinguals from a range of linguistic backgrounds ($M_{\text{age}} = 25.3$, $SD = 4.5$; 13 males). Whilst all participants were considered to be 'young adults' and recruited with corresponding age restrictions, the bilingual group was found to be significantly older than the monolingual group, $t(57.96) = -2.87$, $p = 0.006$, $d = -0.73$. The majority of participants were undergraduate students ($n = 41$), others were postgraduates ($n = 13$) or professionals ($n = 8$), and all but one participant had attended university. All participants gave informed consent prior to testing, had normal or corrected-to-normal vision and did not report to have a history of mental or neurological illness. All bilinguals completed a language history questionnaire adapted from Li, Sepanski and Zhao (2006) with this information summarised in Table 3.1. Based on the self-rated proficiency scores, the bilingual group was highly proficient.

Table 3.1. Bilingual Participants' Language History Information, Experiment 1

Linguistic background	First language	Bulgarian (n = 1)
		Creole (n = 1)
		Dutch (n = 2)
		Farsi (n = 1)
		French (n = 1)
		German (n = 2)
		Hindi (n = 1)
		Hungarian (n = 1)
		Italian (n = 2)
		Lithuanian (n = 1)
		Malayalam (n = 2)
		Polish (n = 7)
		Portuguese (n = 2)
		Romanian (n = 2)
		Sinhalese (n = 1)
		English (n = 4)
Other linguistic background information	Second language	Afrikaans (n = 1)
		English (n = 26)
		Frisian (n = 1)
		Greek (n = 1)
		Gujarati (n = 1)
		Twi (n = 1)
	Third language	English (n = 1)
	Age of first exposure	birth - 6 years (n = 15)
		7 - 12 years (n = 9)
		teenage years (n = 7)
	Time spent in the UK	0 - 5 years (n = 16)
		5 - 10 years (n = 9)
		10+ years (n = 6)
	Switch	rarely (n = 14)
		sometimes (n = 15)
		frequently (n = 2)
Self-rated proficiency (1-6)	Reading	M = 5.1; SD = 0.7
	Writing	M = 4.6; SD = 0.9
	Speaking	M = 4.8; SD = 0.8
	Listening	M = 5.2; SD = 0.7

3.3.6. Participants, Experiment 2

For the second experiment, we recruited 61 participants: 32 English monolinguals and a group of 29 bilinguals. One participant from the bilingual group was excluded because they reported a confidence of 50% on 88% of the trials of the dot discrimination task. Because the confidence marker started at 50% it is likely that this participant simply neglected to provide a confidence judgment for the majority of trials. Including this participant in the non-confidence analyses did not alter the direction or magnitude of any of the effects reported. Therefore, we proceeded to analyse the data provided by a sample of 60 participants, 28 bilinguals and 32 monolinguals. The bilinguals ($M_{\text{age}} = 21.9$, $SD = 4.2$; 6 males) were older than the monolinguals ($M_{\text{age}} = 20.4$, $SD = 0.7$; 7 males), but this difference was only marginally significant $t(28.29)=1.89$, $p = 0.07$, $d = 0.50$).

All of the participants were undergraduate students except for one, who was a postgraduate student. Participants had normal or corrected-to-normal vision and did not report to have a history of mental or neurological illness. All bilinguals completed a language history questionnaire adapted from Li, Sepanski and Zhao (2006). The information from this questionnaire is summarised in Table 3.2.

Table 3.2. Bilingual Participants' Language History Information, Experiment 2

Linguistic background	First language	Bengali (n = 2)
		Cantonese (n = 1)
		Chinese (n = 2)
		English (n = 4)
		French (n = 1)
		German (n = 1)
		Gujarati (n = 1)
		Greek (n = 1)
		Italian (n = 2)
		Korean (n = 1)
		Mandarin (n = 1)
		Nepalese (n = 2)
		Polish (n = 1)
		Portuguese (n = 2)
		Setswana (n = 1)
		Spanish (n = 2)
		Turkish (n = 4)
Other linguistic background information	Second language	English (n = 23)
		Farsi (n = 1)
		French (n = 1)
		Malay (n = 1)
		Punjabi (n = 3)
	Third language	English (n = 2)
		Urdu (n = 1)
	Age of first exposure	birth - 6 years (n = 20)
		7 - 12 years (n = 6)
		teenage years (n = 3)
	Time spent in the UK	0 - 5 years (n = 13)
		5 - 10 years (n = 2)
		10+ years (n = 14)
	Switch	rarely (n = 13)
		sometimes (n = 14)
		frequently (n = 2)
Self-rated proficiency (1-6)	Reading	M = 5.0; SD = 0.9
	Writing	M = 4.7; SD = 1.1
	Speaking	M = 4.8; SD = 1.1
	Listening	M = 5; SD = 1.0

3.3.7. Hierarchical Models

Note that all predictors entered into the hierarchical models are z-scored on the participant level, and that response time was log-transformed prior to being z-scored to make RT distributions approximately normal. All models reported in this chapter allowed for random intercepts and random slopes at the participant level.

3.4. Results

3.4.1. Control Measures, Experiment 1

An analysis of the control measures revealed that both groups performed comparably on measures of working memory, $t(56.17)=1.67, p=.10, d=0.42$ and nonverbal reasoning, $t(59.62)=0.74, p=.46, d=0.19$. Means and standard deviations are reported in Table 3.3. Therefore, any differences found in metacognitive abilities are unlikely to be attributable to group differences in general cognitive functioning.

Table 3.3. Descriptive Statistics for Control measures, Experiment 1

	Monolinguals		Bilinguals	
	M	SD	M	SD
Working Memory (maximum score: 30)	17.97	4.85	16.03	3.73
Nonverbal Reasoning (maximum score: 20)	9.94	1.65	10.26	1.79

3.4.2. Control Measures, Experiment 2

Both groups performed comparably on measures of working memory, $t(56.49)=-0.91, p=0.37, d=0.23$ and nonverbal reasoning, $t(57.20)=0.98, p=0.33, d=0.25$, indicating that the groups were matched on general cognitive functioning (see Table 3.4.).

Table 3.4. Descriptive Statistics for Control Measures, Experiment 2

	Monolinguals		Bilinguals	
	M	SD	M	SD
Working Memory (maximum score: 30)	15.66	3.55	14.93	2.62
Nonverbal Reasoning (maximum score: 12)	8.66	2.47	9.21	1.91

3.4.3. First Order Performance, Model Free Analyses

We compared the bilinguals' and monolinguals' performance with regards to their first order accuracy (percentage of correct responses), the difficulty of the trials (dot difference) and the response time of the choices (ms) for both the first and second experiment (see Figure 3.3. below). Monolinguals and bilinguals had similar proportions of correct responses in the first experiment, $M_{\text{monolingual}}=70.1\%(1.1\%)$, $M_{\text{bilingual}}=70.1\%(1.2\%)$, $t(58.73)=0.66$, $p=0.51$, $d=0.17$. However, for the second experiment monolinguals had a slightly higher proportion of correct responses $M_{\text{monolingual}}=71.3\%(1.0\%)$, $M_{\text{bilingual}}=70.3\%(1.4\%)$, $t(56.78)=3.44$, $p=0.001$, $d=0.88$. These analyses suggest that the staircase procedure worked, as all groups are close to 70% correct with small standard deviations. The mean dot difference was similar for both groups in Experiment 1 ($M_{\text{monolingual}}=4.34$ (1.08), $M_{\text{bilingual}}=4.46$ (1.03), $t(59.91)=1.25$, $p=0.27$, $d=0.29$) and Experiment 2 ($M_{\text{monolingual}}=5.85$ (1.80), $M_{\text{bilingual}}=5.46$ (1.22), $t(54.82)=1.02$, $p=0.31$, $d=0.26$), meaning that the task was equally difficult for both groups. Note that the dot difference was higher for the second experiment than the first experiment, reflecting that performance dropped under speed stress $t(103.82) = 4.90$, $p < 10^{-5}$, $d = 0.89$). Finally, monolinguals took longer to respond than bilinguals in Experiment 1, when response times were unconstrained ($M_{\text{monolingual}}=3359$ (1474), $M_{\text{bilingual}}=2679$ (922), $t(50.38)=2.18$, $p=0.03$, $d=0.55$) but response times were similar for both groups for Experiment 2 ($M_{\text{monolingual}}=856$ (150), $M_{\text{bilingual}}=898$ (95), $t(52.91)=1.29$, $p=0.20$, $d=0.33$). Unsurprisingly, responses were on average faster in the second experiment compared to the first ($t(62.29) = 13.24$, $p < 10^{-10}$, $d = 2.38$). In Experiment 2, trials were excluded if participants failed to respond within 1.5 seconds. Both language groups lost a similar proportion of trials due to slow responding (3% for monolinguals and 4% for bilinguals; $t(52.82) = 1.85$, $p = .07$, $d = 0.48$).

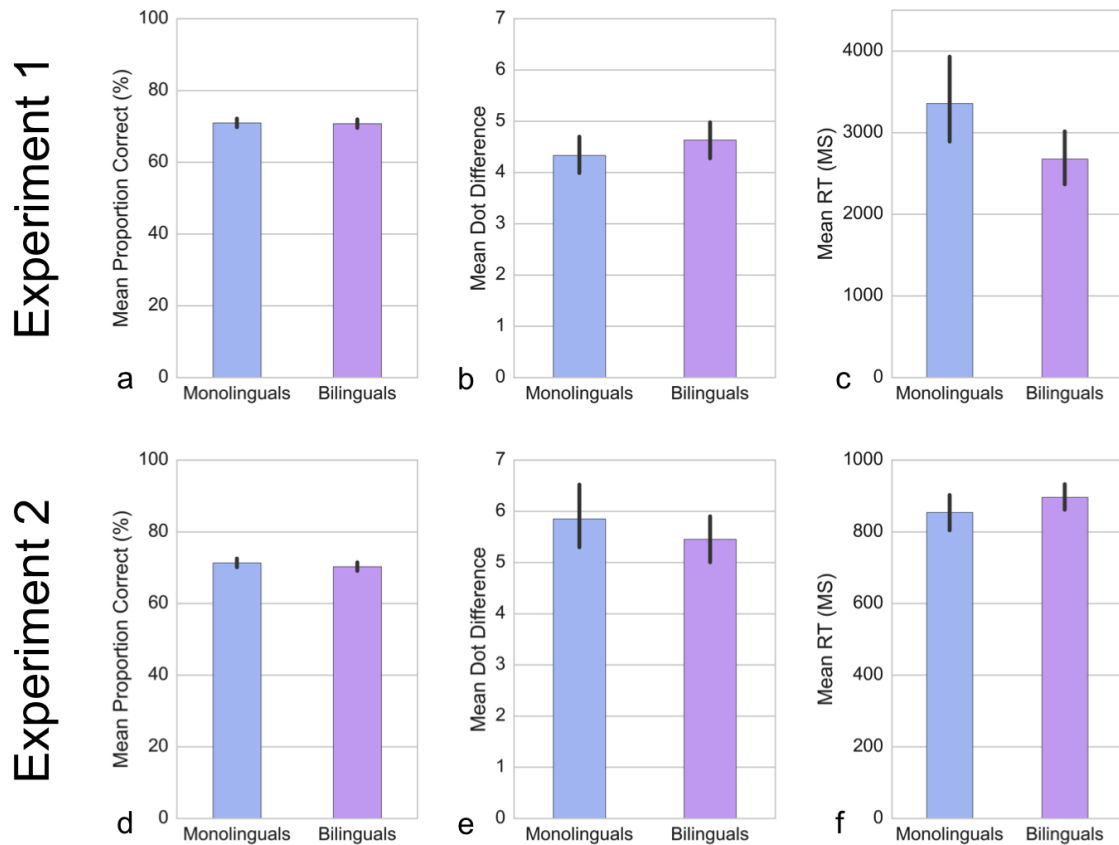


Figure 3.3. Comparing First Order Performance of Monolingual and Bilingual Participants

Bar graphs comparing the mean values of monolinguals and bilinguals for proportion correct responses (**a & d**) dot difference (**b & e**) and response time (**c & f**, note the difference in y scaling), for the first and second experiment, respectively. Mean proportion correct is close to the 70% target for both groups in both experiments, but it is significantly higher for monolinguals than bilinguals in Experiment 2. Mean dot difference was similar for both groups in both experiments but Experiment 2 had higher dot difference than Experiment 1, signifying that the response time constraint of the second experiment reduced participants' ability to discriminate between the stimuli. Monolinguals responded significantly slower in Experiment 1 response times were similar for both groups in Experiment 2 when a response time limit of 1.5 seconds was enforced.

To test how response time and dot difference influenced first order accuracy I ran hierarchical logistic regression models. The models predicted the probability of a choice being correct from z-scored dot difference and log-transformed response time and interaction terms between these predictors and the language group of the respondent (dummy coded as 1 for the monolingual group and 0 for the bilingual group). Dot difference predicted accuracy for Experiment 1

($z=15.34$, $p<10^{-10}$) and Experiment 2 ($z=13.55$, $p<10^{-10}$) and there was no difference in the slopes between the language groups (Experiment 1: $z=0.02$, $p=.98$; Experiment 2: $z=-0.03$, $p=.98$; See Figure 3.4.).

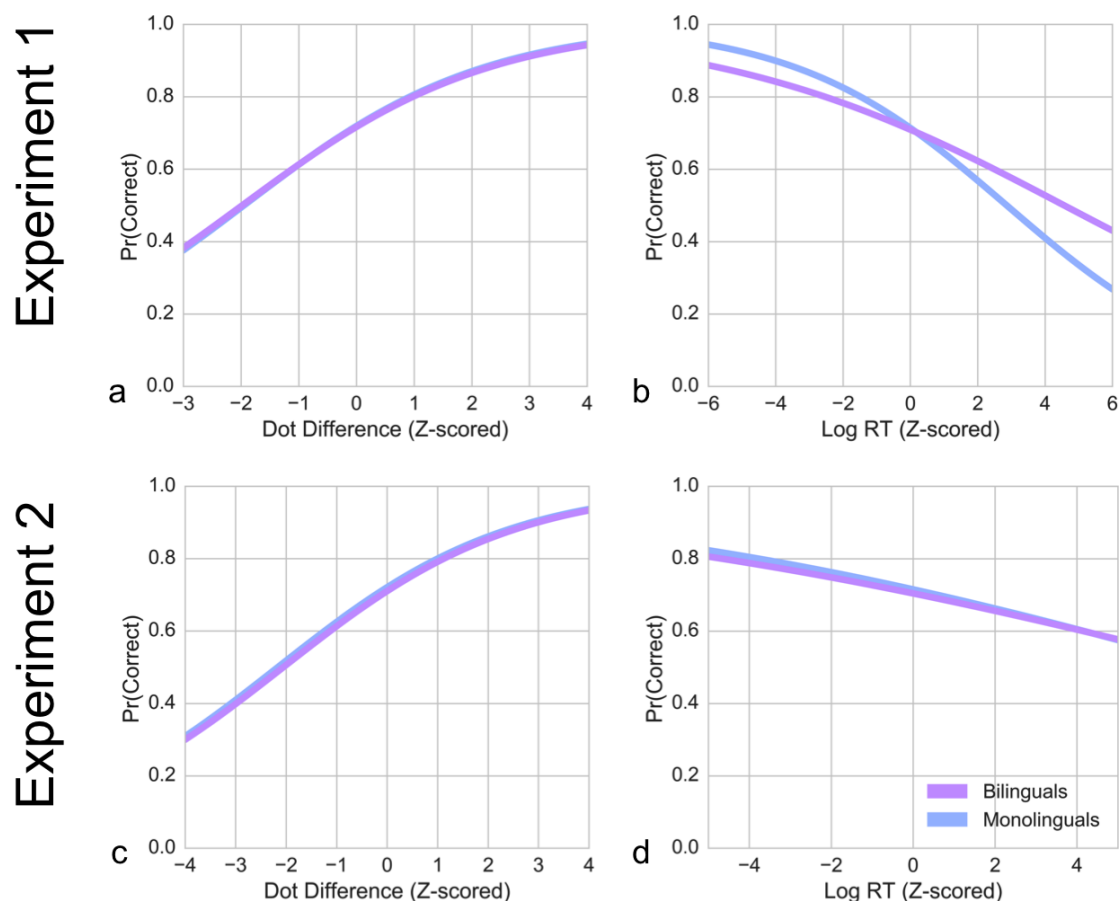


Figure 3.4. Predicting First Order Accuracy From Dot difference and Response Time

Logistic regression curves predicting accuracy (probability of being correct) from dot difference (a & c) and response time (b & d) for Experiment 1 and 2, respectively. Dot difference predicted accuracy for both experiments. There was no difference between the language groups with regards to the strength to this relationship. Response time predicted accuracy for both experiments. This relationship were stronger for the monolingual group in Experiment 1 (where monolinguals had slower response times than bilinguals) but not in Experiment 2 (where both groups had similar response times).

Response time predicted accuracy for both experiments (Experiment 1: $z=-5.95$, $p<10^{-8}$; Experiment 2: $z=-2.02$, $p=.04$), the slope is steeper for the monolingual group in Experiment 1 ($z=-2.54$, $p=.01$), with no difference between the groups in Experiment 2 ($z=0.03$, $p=.98$). In sum, first order performance appear to be comparable between the groups with the notable

exception of the response times in the first experiment and the proportion correct in the second experiment. To test whether these differences reflected different rates of evidence accumulation in the first order task I ran a hierarchical drift diffusion model.

3.4.4. First Order Performance, DDM

All data from the main trials was entered into the DDM, using the HDDM package for Python (see Chapter 2 for more information). Drift rate was determined by dot difference, a sensitivity parameter and an intercept. Boundary separation, non-response time and the drift rate intercepts were allowed to vary between groups. Because of the hierarchical nature of the model, each participant had an individual parameter estimate for boundary separation, non-decision time and drift rate drawn from a distribution determined by which group they belonged to. Monolinguals had a greater boundary separation than bilinguals with a 94% probability in Experiment 1, corresponding to posterior odds of about 16:1. This means that monolinguals emphasised accuracy over speed when they responded (see Figure 3.5.). These parameter estimates fit the raw data where monolinguals were slower but had a (nonsignificantly) lower dot difference than their bilingual peers. In the second experiment this pattern was reversed: Bilinguals had a greater boundary separation than monolinguals with an 89% probability (posterior odds=8:1).

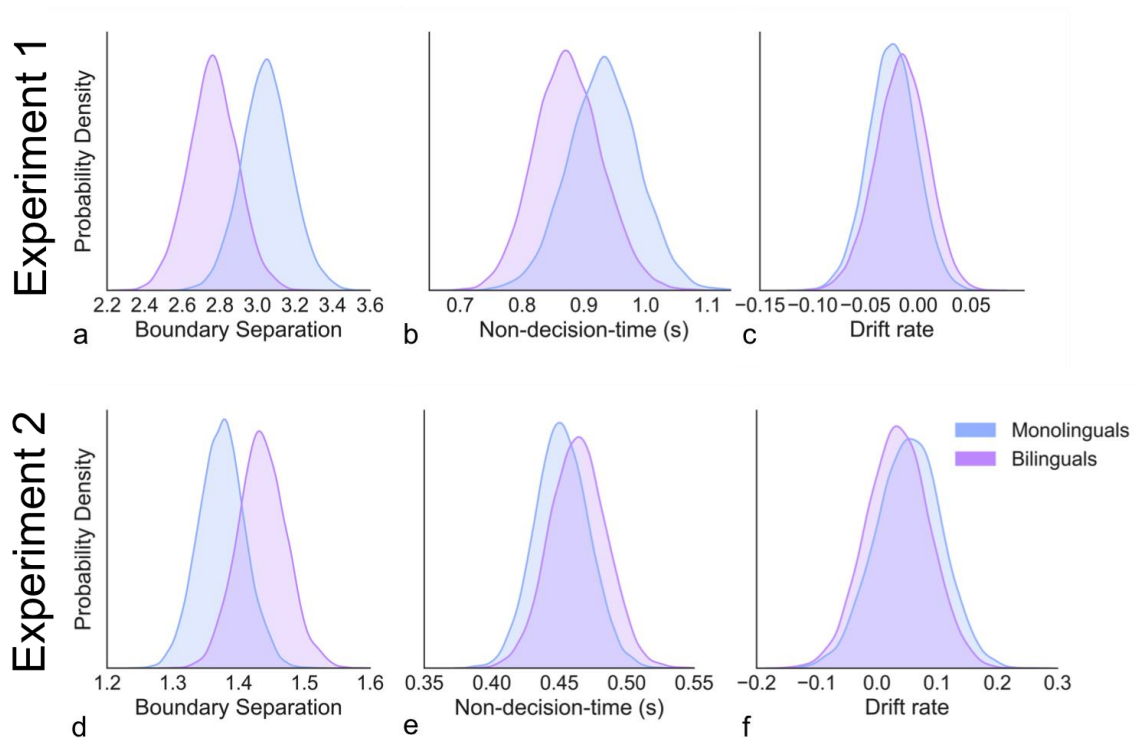


Figure 3.5. Parameter Estimates from the DDM Model

Probability densities for boundary separation (a & d), non-decision-time (b & e), drift rate (c & f)

f). There were no systematic differences between the groups in boundary separation, non-decision times or drift rates across experiments.

Non-decision times were longer for monolinguals than bilinguals in Experiment 1 with a 78% probability (posterior odds=3:1). The direction of this difference was reversed in Experiment 2 where bilinguals had longer non-decision times with 67% probability (posterior odds=2:1), but it is also worth noting that the two distributions show 77% overlap. Finally, drift rates show 85% overlap in Experiment 1 (posterior odds=2:1) and 86% overlap in Experiment 2 (posterior odds=2:1). To summarise, Monolinguals had greater boundary separation in Experiment 1 but bilinguals had greater boundary separation in Experiment 2. Bilinguals have longer non-decision times in the first experiment but bilinguals appear to have slightly longer non-decision time in Experiment 2. Drift rates appear to be similar across the groups in both experiments. Together, these findings suggest that there are no systematic differences between monolinguals and bilinguals with regards to first order response strategies or first order performance.

3.4.5. Second Order Performance

To explore how monolinguals and bilinguals differed with regards to raw confidence scores, ranging from 0 for the lowest possible confidence judgement to 100 for the highest possible confidence judgement. I ran a hierarchical linear regression model that predicted confidence from accuracy, language group and an interaction between confidence and language group, while allowing the intercepts and the effect of accuracy to differ by participant. Confidence was on average higher for correct trials than for error trials in both Experiment 1 ($t(60.30)=7.61, p>10^{-9}$) and in Experiment 2 ($t(57.96)=7.9, p>10^{-10}$). Despite both groups having a similar first order performance, bilinguals reported feeling more confident than monolinguals in Experiment 1 ($t(60.04)=4.05, p=.001$) and Experiment 2 ($t(57.87)=1.78, p=.08$) when first-order accuracy was controlled for. This difference appears to be particularly pronounced for error trials (see Figure 3.6.), but this trend is not significant in either Experiment 1 ($t(60.41)=1.92, p=.06$) or Experiment 2 ($t(58.11)=1.67, p=.10$). In other words, it seems as if monolingual confidence judgements are more sensitive to first-order accuracy than bilingual confidence judgements but these differences are not significant in these rudimentary analyses.

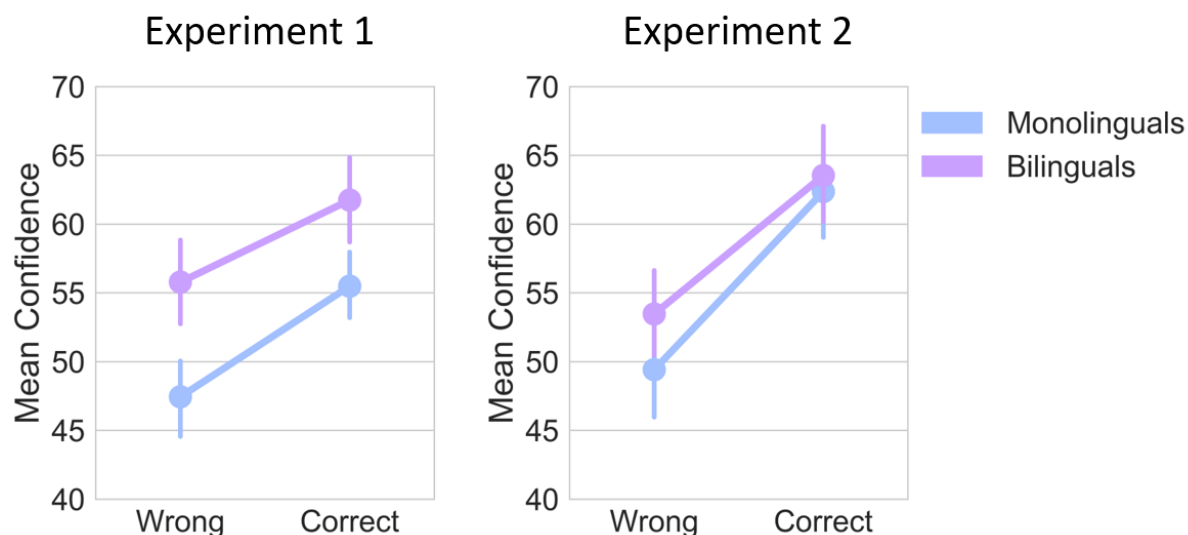


Figure 3.6. Raw confidence judgements by monolinguals and bilinguals, as a function of the accuracy of the response. Monolinguals report being more confident than bilinguals in Experiment 1 and Experiment 2. Both groups report being more confident when they are correct than when they are wrong, and in both experiments there is a trend that monolingual confidence judgements are more sensitive to first-order accuracy than bilingual confidence judgements, but this trend is insignificant.

In order to further examine this potential difference between monolinguals and bilinguals in metacognitive processing I needed a more sensitive measure. Signal detection theory offers a way to estimate metacognitive efficiency, which captures how good a person's metacognitive accuracy is when controlling for their first order accuracy (See Chapter 2). To estimate metacognitive efficiency I used the *Mratio*, a ratio between metacognitive sensitivity, captured by $\text{meta-}d'$, and first-order sensitivity, captured by d' . This measure is useful because it accounts for variation in first-order sensitivity, first-order response thresholds and second order response thresholds, something that correlational measures, such as regression coefficients fail to do (Barret et al., 2013; Galvin et al., 2003). I used a hierarchical Bayesian estimation method to fit the *Mratios*. The Bayesian estimation method has the advantage that it explicitly represents the group level-parameters as probability distributions that accounts for the uncertainty at each level of estimation (Fleming, 2017; See Chapter 2). The confidence data was binned into participantwise tertiles prior to the *Mratio* estimation. I fitted each language group from each experiment separately, to prevent shrinkage from hiding any true group difference. I did not run a t-test on the individual *Mratio* estimates, as each *Mratio* estimate was constrained by the hyperparameter of that group, thus violating the independence assumption. Instead I operated on the posteriors of the group means directly. For each experiment I subtracted the posterior of the

mean Mratio of the bilingual group from the posterior of the mean for the monolingual group, creating new posteriors for the Mratio difference, quantifying how much more (or less) metacognitively efficient monolinguals were relative to bilinguals. 95% of the probability mass of the mean Mratio difference was greater than 0 meaning that monolinguals were more metacognitively efficient than bilinguals (posterior odds= 21:1). The second experiment replicated this effect as 95% of the probability mass was greater than 0, giving 20:1 odds that monolinguals were more metacognitively efficient than bilinguals. Because both experiments had a probability distribution over the parameter of interest (the difference in metacognitive efficiency between the language groups) I wanted to use Bayes Rule to combine the information from both posteriors to get a more accurate estimate of the population difference. I am not aware of a method to update from two empirical posteriors. However, because both empirical posteriors were approximately normal (see Figure 3.6.), I could first estimate analytic normal distributions that fitted the numeric posteriors well, by taking the mean and standard deviation of each numeric posterior distribution, and then analytically combine these two analytic normal distribution into a new hyper posterior. Because the normal distribution is conjugate with itself, finding mu and sigma of the new distribution was analytically tractable. The precision of the hyper posterior was given by adding the precisions of the two experimental posteriors. Because the precision is the inverse of the variance, it was straightforward to compute the standard deviation of the hyper posterior. The mean of the hyper posterior was the precision-weighted combination of the means of the two empirical posteriors. The hyper posterior, representing the best estimate of the population difference in metacognitive efficiency between monolinguals and bilinguals had 1% of its probability mass below 0 (posterior odds=97:1). In other words, the combined information of these two experiments provide strong evidence that monolinguals are more metacognitively efficient than bilinguals, at least in the context of visual discrimination.

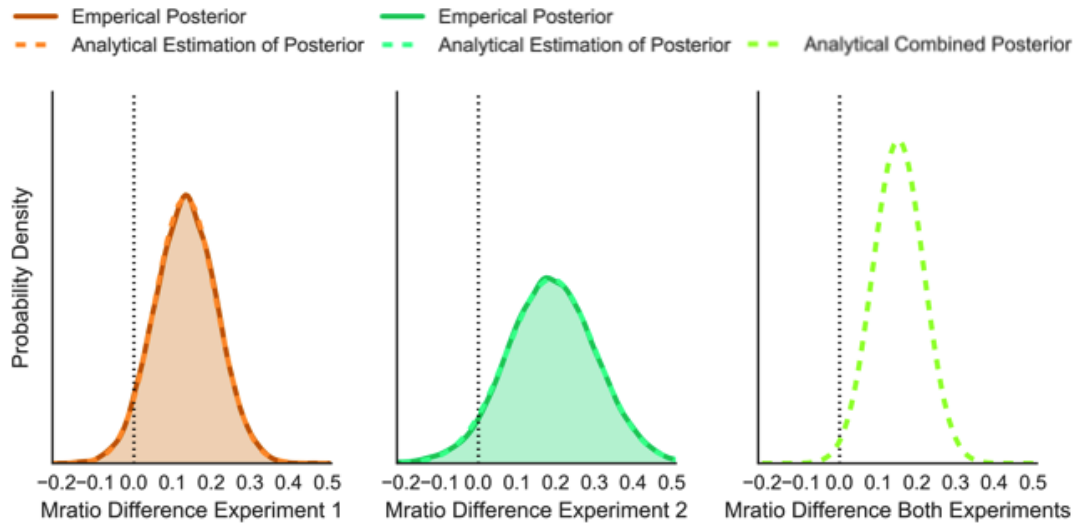


Figure 3.7. Difference in Metacognitive Efficiency Between Monolinguals and Bilinguals

The estimated mean difference in Mratios between monolinguals and bilinguals based on Experiment 1, Experiment 2 and the combined data from both experiments. The Mratio is a signal-detection-theoretic measure that captures metacognitive accuracy when first order accuracy and metacognitive response-thresholds are controlled for. Mratios were computed using a numeric Bayesian estimation method (the filled curves). The combined posterior was computed analytically, by fitting a normal distribution to Experiment 1 and Experiment 2 (the dashed lines) and using Bayesian updating to combine the information in both. The x-axes show the mean monolingual Mratio – the mean bilingual Mratio. Mratios were higher for the monolingual group than the bilingual group in Experiment 1 and Experiment 2.

Because monolinguals had both higher Mratios and slower response times in Experiment 1, I wanted to test whether mean response times were associated with Mratios (Figure 3.8.) In both experiments, response times did not predict Mratios when group affiliation was accounted for (Experiment 1: $t(59)=-0.46$, $p=.65$; Experiment 2: $t(57)=-0.52$, $p=.61$). Finally, I predicted Mratios from self-reported second language proficiency among the bilinguals. In the second experiment self-reported proficiency predicted metacognitive accuracy ($t(26)=2.41$, $p=.02$), but this was not the case for the first experiment ($t(29)=0.06$, $p=.95$).

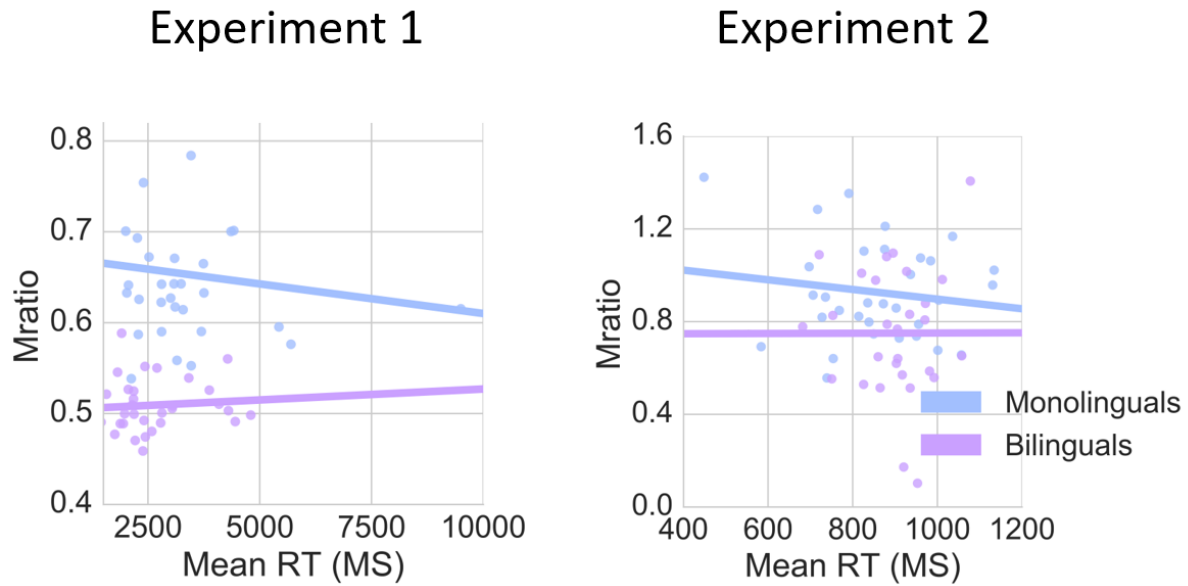


Figure 3.8. Mratios as a function of mean response time and group affiliation

Mratios (measures of metacognitive efficiency) are not associated with the mean response times of participants in either experiment when their group affiliation is accounted for.

To further investigate potential causes for the difference in Mratio, I explored how RT and dot difference influenced confidence for each group (see Figure 3.9., a and b) by predicting z-scored confidence from z-scored RT and z-scored dot difference. I allowed the slopes of RT and dot difference to vary by participant, the model contained no intercept parameters because the purpose of the model was to explore the relationship between the variables. Finally, the model contained interaction terms that allowed the effect of RT and dot difference on choice to vary by language group. RT negatively predicted confidence in Experiment 1 ($t(60)=-10.37, p<10^{-10}$) and Experiment 2 ($t(58.03)=-6.60, p<10^{-7}$). There was no statistically significant difference between monolinguals and bilinguals with regards to the strength of the relationship between RT and confidence either in Experiment 1 ($t(60)=1.47, p=.15$) or Experiment 2 ($t(58.16)=1.12, p=0.27$). Dot difference was a positive predictor for both experiments (Experiment 1: $t(60)=4.92, p<10^{-5}$; Experiment 2: $t(57.54)=6.08, p<10^{-6}$) and the strength of this relationship was similar for both groups (Experiment 1: $t(57.51)=-0.25, p=.80$; Experiment 2: $t(60)=0.25, p=.81$).

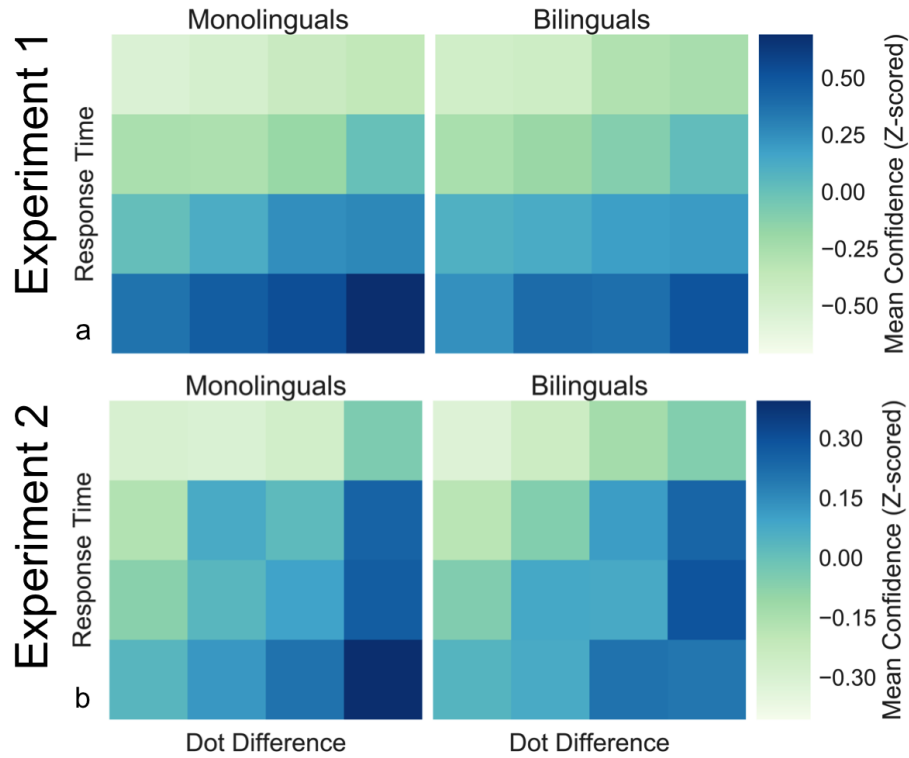


Figure 3.9. Predictors of Confidence

(a & b) shows how response time and dot difference influence confidence for the monolingual and bilingual groups. Response time and dot difference are binned into participantwise quantiles and the mean z-scored confidence judgement of each cell is presented. Response time and dot difference seems to influence confidence in a linear and additive fashion.

To further understand what influenced confidence in the different language conditions I split the data into correct trials and error trials and examined how RT and dot difference influenced confidence in the two language groups (Figure 3.10.). I modelled this by running separate models for correct and error trials for each experiment. Apart from predicting confidence for errors and correct choices separately the models were equivalent to the hierarchical regression models described in the previous section. In Experiment 1 RT negatively predicted confidence for both correct ($t(59.89)=-9.64, p<10^{-10}$) and incorrect trials ($t(58.77)=-7.37, p<10^{-9}$). This effect was similar for both language groups (correct: $t(59.59)=1.37, p=.18$; error: $t(59.16)=0.91, p=.37$). In Experiment 2 slower responses were also associated with lower confidence for correct ($t(57.01)=-7.10, p<10^{-10}$) and incorrect trials ($t(52.31)=-2.80, p=.007$). The influence of RT on confidence was marginally weaker for the bilingual group than the monolingual group ($t(57.91)=1.67, p=.10$) for the correct trials. Note that this effect was in the same direction in Experiment 1 but insignificant. There was no difference in the influence of confidence on RT for the error trials ($t(53.07)=0.01, p=.99$).

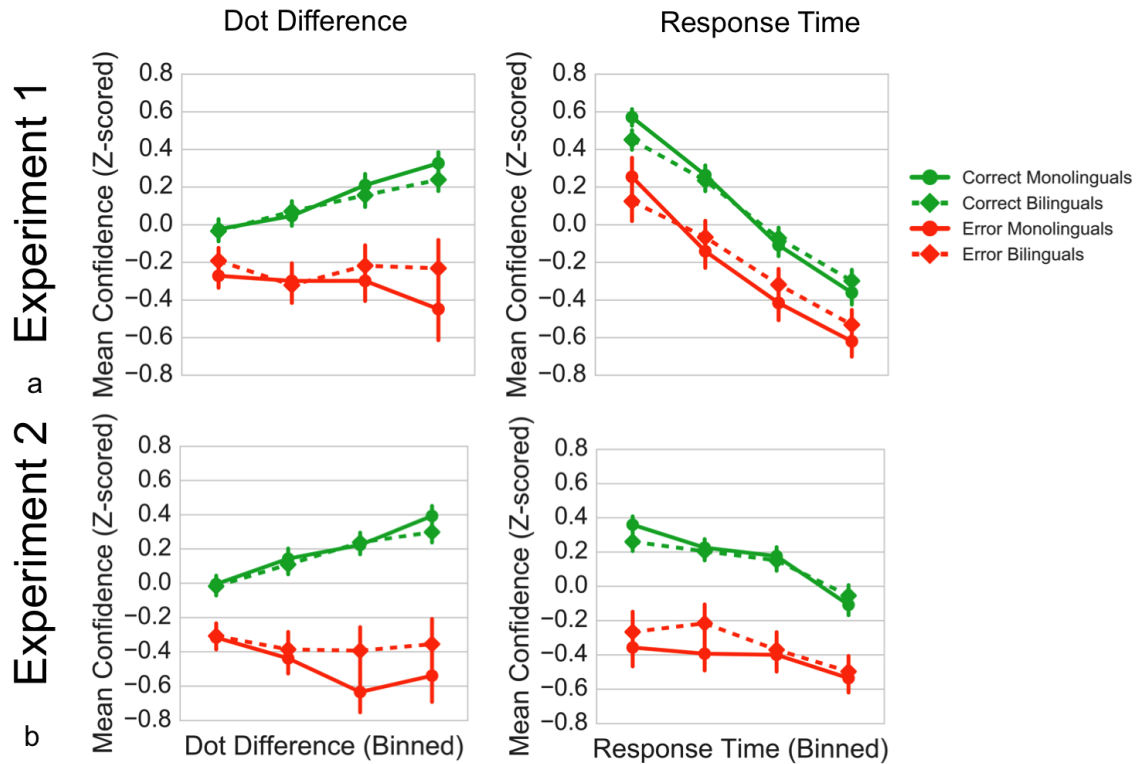


Figure 3.10. Predictors of Confidence by Accuracy

(a & b) shows how response time and dot difference influence confidence for the monolingual and bilingual groups as a function of accuracy. Dot Difference is binned into quartiles from lowest to highest. Response times are quartile binned from fastest to slowest. Easier trials (trials with higher dot difference) are associated with higher confidence for correct trials, but lower confidence for error trials. This pattern holds for both experiments and both groups. Slower response times are associated with lower confidence for both language groups for correct and error trials in Experiment 1. The same pattern is present but weaker in Experiment 2, when response times were constrained.

Dot difference positively predict confidence for correct trials in Experiment 1 ($t(58.95)=5.19$, $p<10^{-5}$) and Experiment 2 ($t(58.69)=6.26$, $p<10^{-7}$). The strength of this effect is comparable across groups (Experiment 1 $t(59.01)=0.49$, $p=.61$; Experiment 2: $t(58.22)=0.33$, $p=.74$). Dot difference does not predict confidence for error trials in Experiment 1 ($t(69.85)=-0.83$, $p=.41$) but has a marginal negative effect of confidence in Experiment 2 ($t(64.04)=-1.86$, $p=.07$). The relationship between dot difference and confidence during error trials is not mediated by language group in either experiment (Experiment 1: $t(64.82)=-0.12$, $p=.91$; Experiment 2: $t(55.62)=-0.92$, $p=.35$).

3.4.6. Exploring Potential Non-linear Relationships

The previous sections point to an interesting conundrum, the monolinguals are more metacognitively accurate than their bilingual peers their confidence judgments do not appear to be more sensitive to stimulus strength (captured by dot difference) or response time. This suggests that there is some third variable that is diagnostic of accuracy that differentially affects confidence between the groups. However, it is also possible that the reason I have failed to find a difference in sensitivity to response time and stimulus strength is because I have only tested linear models. To examine whether the linearity assumption hid group differences I used a cubic smoothing spline on the pooled z-scored data to estimate the non-linear interaction between stimulus strength and response time in predicting accuracy and confidence (Experiment 1: Figure 3.11, Experiment 2: Figure 3.12.). A smoothing spline is a method for fitting a non-linear curve (or in the case of two predictors a non-linear plane) to a set of observations, by satisfying a cost function with two components, minimising the discrepancy between the predicted values and the observed outcomes and keeping the curve (plane) smooth, the relative weight of minimising errors relative to maximising smoothness is determined by a smoothing parameter λ (Friedman et al., 2001). λ Was here determined by generalised leave one out cross validation, which means that the smoothing parameter that led to the lowest out of sample prediction error was selected. Consequently I got a near optimal approximation of how stimulus strength and response time related to accuracy and confidence in these two experiments. I then fitted the predicted values (probability correct and z-scored confidence) to an equidistant 41x41 grid ranging between -2 and +2 z-scored response time units on one axis and -2 and +2 z-scored dot difference units on the other axis. My aim was to see if there were any systematic discrepancies between how RT and stimulus strength related to confidence on the one hand and accuracy on the other that could account for the observed differences in metacognitive accuracy between monolinguals and bilinguals.

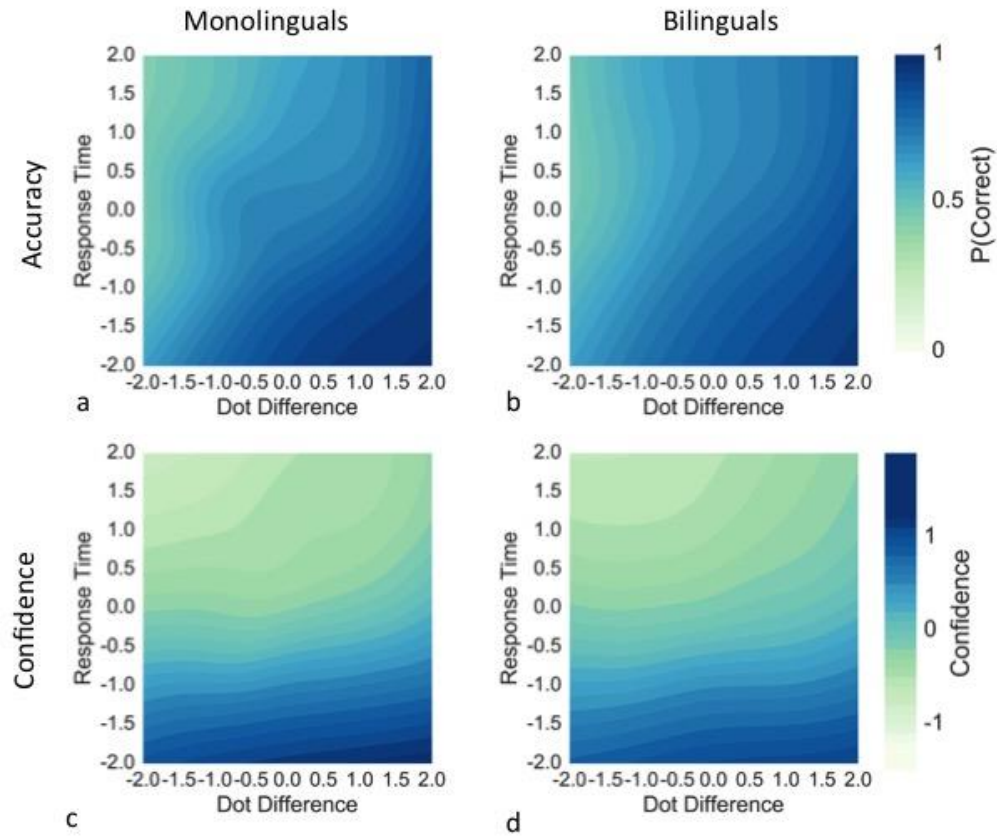


Figure 3.11. Map of the Non-linear Relationship between Response Time and Dot Difference in Predicting Accuracy and Confidence, Experiment 1

(a & b) shows how response time and dot difference influence accuracy for the monolingual and bilingual groups. (c & d) show confidence as a function of response time and accuracy for the monolingual and bilingual groups. Values for accuracy and confidence were simulated for a 41x41 response time and dot difference grid, based on the best estimate of a smoothing spline model fitted to the whole data. Response time, dot difference and confidence were z-scored. Dot difference appear to be predict accuracy more strongly than response time but dot difference and response time appear to predict confidence roughly to the same extent.

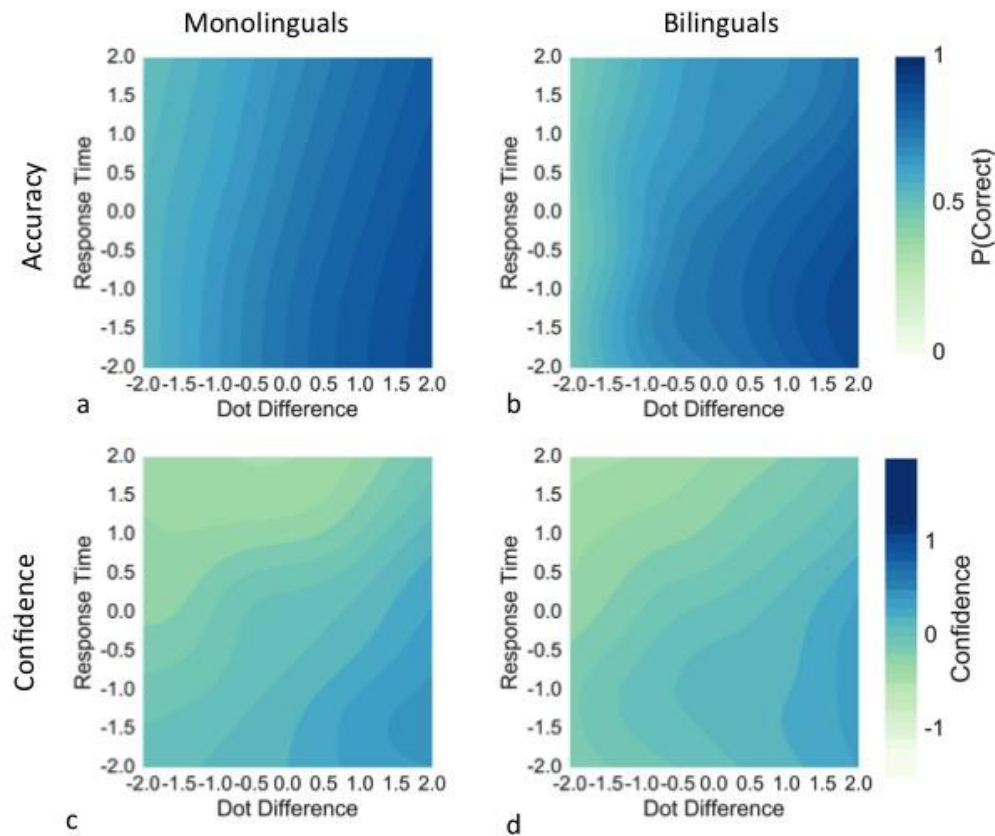


Figure 3.12. Map of the Non-linear Relationship Between Response Time and Dot Difference in Predicting Accuracy and Confidence, Experiment 2

(a & b) shows how response time and dot difference influence accuracy for the monolingual and bilingual groups. (c & d) show confidence as a function of response time and accuracy for the monolingual and bilingual groups. Values for accuracy and confidence were simulated for a 41x41 response time and dot difference grid, based on the best estimate of a smoothing spline model fitted to the whole data. Response time, dot difference and confidence were z-scored. Dot difference appear to be predict accuracy more strongly than response time, but dot difference and response time appear to predict confidence roughly to the same extent.

It is noteworthy that while accuracy seems to be equally predicted by dot difference and response time in Experiment 1, confidence seems to be dominated by response time. Perhaps because participants have a more reliable internal estimate of their response times than they do of the dot difference? (After all, if they had a perfect internal representation of the dot difference they would also have perfect performance). On the second task, in which I had imposed a time limit on responding, z-scored response time is much less diagnostic of accuracy, in line with the

results from the logistic models reported above. Interestingly, participants seem to have adapted their response criterion for confidence, taking dot difference more into account, but as in experiment 1 they overweight the influence of response time in their confidence judgments. All in all, these results closely match the results of the linear models, suggesting that the linearity assumption has not unduly biased the findings. It is also clear from the graphs that the overall patterns are similar for monolinguals and bilinguals in both groups, so the monolingual metacognitive advantage cannot be explained by non-linear interactions between response time, difficulty and confidence.

3.4.7. Control Analyses

Because bilinguals were significantly older than monolinguals in Experiment 1, and marginally older in Experiment 2, I controlled for age when predicting Mratios from language group, using an ordinary least squares regressions. Age did not significantly predict Mratios for Experiment 1 ($t(59)=0.56$, $p=.57$) nor Experiment 2 ($t(57)=1.13$, $p=.27$). Neither did it modify the strength of the relationship between group and Mratio for either experiment (Experiment 1: $t(59)=8.74$, $p<10^{-11}$; Experiment 2: $t(57)=2.96$, $p<.005$). Similarly, because working memory scores were marginally different between the groups in Experiment 1, I controlled for working memory when looking at differences in Mratio by language group. Working memory did not predict Mratios ($t(59)=1.49$, $p=.14$) and did not influence the effect of Mratio by language group ($t(59)=8.74$, $p<10^{-11}$).

3.5. Discussion

In two experiments I compared the metacognitive ability of young, healthy monolinguals and bilinguals in a perceptual 2AFC paradigm, in which participants tried to determine which of two circles contained the most dots and subsequently made a confidence judgment about their accuracy. The first experiment had unconstrained response times, whereas the second experiment enforced a 1.5-second response threshold on first order responses. The aim of this constraint was to keep first order performance similar across as many dimensions as possible.

Participants' response times in Experiment 2 were significantly faster in both groups relative to Experiment 1. This increase in response speed came at a cost to performance, illustrating the well-established trade-off between speed and accuracy (Forstmann et al., 2016; Vickers & Packer, 1982). Given the staircase procedure implemented to keep the error-rate constant, this performance difference resulted in easier trials (i.e., a greater dot difference) in Experiment 2 compared to Experiment 1. When people did not respond within 1.5 seconds in Experiment 2

the trial was discarded. The monolingual and bilingual groups missed a similar proportion of trials because of responding too slowly (4% or less), and the two groups did not differ significantly with regards to average trial difficulty.

Despite the staircase procedure constraining the variance in difficulty, dot difference was a strong predictor of accuracy in both experiments, and this effect was equally strong for both groups in both experiments. In the first experiment z-scored RT was a stronger predictor of accuracy for the monolingual group. However, this effect might just be because there was a greater range of response times in the monolingual group, so the z-scored response times had a higher resolution. For the second experiment, where the range of response times were similar between the groups there was no difference in the predictiveness of z-scored RT. Fitting a hierarchical DDM model to the data showed that the rate of evidence accumulation was similar across language groups for both experiments, and while the boundary separation appeared to be farther and non-decision times seemed to be longer for the monolinguals in the first experiment, these patterns reversed for the second experiment. This reversal suggests that there are no systematic differences in response criterion and non-decision times between the groups. In summary, first order performance was very similar across the two groups, and there is no reason to believe that the evidence accumulation leading up to the first order decision was systematically different for bilinguals and monolinguals.

Bilinguals were on average more confident than monolinguals in both experiments, but because the participants were instructed to report their relative confidence rather than their probability of being correct it is impossible to say which group was better calibrated. With regards to metacognitive efficiency, the monolingual group had significantly higher M_{ratios} than the bilingual group, in both experiments. These differences could not be explained in terms of differences in non-verbal reasoning, working memory or age. Because the monolinguals responded more slowly in Experiment 1 and also had higher M_{ratios} , I investigated whether slower mean response times were associated with higher M_{ratios} , the hypothesis being that slower response times would allow more time for meta-cognitive processing (see Baranski & Petrusic, 1998, 2001). In both experiments mean response times did not seem to be associated with metacognitive efficiency when the difference between the language groups was accounted for.

To examine why monolinguals were more metacognitively sensitive than bilinguals I predicted confidence from difficulty and response time for both groups, splitting the data into correct and error trials. There was no significant differences between the groups, either for correct or error

trials in either experiment. Slower responses were associated with lower confidence both for correct and incorrect trials for both experiments, but this effect was weaker in Experiment 2, probably because the range of response times was artificially constrained by the experimental design. Dot difference positively predicted confidence judgments for both experiments for correct trials, but not for error trials (and there was a marginally significant negative relationship in Experiment 2). These findings support the idea that both response time and stimulus strength influence confidence judgments, but does so in different ways (see Kiani, Corthell, & Shadlen, 2014). Finally, I ran a smoothing spline on dot difference and response time, predicting accuracy and confidence, to test if the previous linear models had missed any non-linear differences between the monolingual and bilingual groups. Visual inspection suggests that is not the case.

Since I expected to find a metacognitive advantage for bilinguals, it might seem like these results are at odds with the broader literature on bilingualism. This is to some extent true, but should not be overstated. As discussed in the introduction, the executive function advantage in bilinguals has not been consistently replicated across studies, and the exact nature of this advantage is still controversial (de Bruin et al., 2015; Paap et al., 2015).

It is still unclear what is causing the bilingual disadvantage in metacognitive sensitivity. Self-reported second language ability did predict metacognitive efficiency in Experiment 2, but not in Experiment 1. The fact that this relationship is not reliable across studies, together with the small sample size, suggests that it is probably a false positive, especially considering that self-reported language proficiency is quite a noisy measure of actual competence (MacIntyre, Noels, & Clément, 1997). In order to better understand when during language learning this effect manifests one of my colleagues is running a follow-up study that examines metacognition in monolingual and bilingual school children of various ages. Though the actual cause of the observed difference is an open question, the current study has helped rule out some possible explanations. Specifically, it appears that the difference in metacognitive accuracy cannot be accounted for by differences in response speed, speed of first order evidence accumulation or how response time and difficulty influences confidence. This third point is of particular interest because historically response time and stimulus strength have been presumed to be the two primary drivers of confidence judgments (Kiani et al., 2014; Van Den Berg et al., 2016; Vickers & Packer, 1982). However, that they can't explain the difference in metacognitive performance between monolinguals and bilinguals suggests the presence of some third variable that is diagnostic of accuracy and influences confidence, and acts differently on monolinguals and bilinguals. Hunting for this third variable would be a worthwhile research project, and studies that have been published after this research was conducted provide some interesting avenues.

First, Siedlecka and colleagues (2016) found that the accuracy of confidence judgments differed based on whether they were recorded before or after the choice. Fleming and Daw (2017) have argued that this is because the choice itself provides information to the metacognitive system, provided that the evidence streams for choice and confidence are at least partially decoupled. This means that monolingual confidence judgments may be more accurate than bilingual confidence judgments because monolinguals are better at integrating the information from the choice itself into the confidence judgment. This could be tested by a task where participants have to give confidence ratings either before or after they make a choice. If bilingual and monolingual metacognitive accuracy is similar for confidence judgments before they make a choice, but different for confidence judgments after they make a choice that would support this explanation. This idea will be discussed further in the general discussion. Alternatively, I might fail to find a difference in sensitivity to stimulus strength because the measure of stimulus strength is too noisy. Stimulus strength was here conceptualised as dot difference, but it is probable that other aspects of the stimuli, such as dot clustering also influence difficulty (it is harder to estimate the amounts of dots in a circle when they are overlapping). Additionally the current design did not allow me to track or control how the participants sampled evidence. Random dot kinematograms would have solved this problem as each “movement” can be treated as an independent sample from a normal distribution, and average motion strength and variation can both be controlled and recorded.

Another new avenue of inquiry would be how these differences in metacognitive sensitivity translate to higher-order decision making. Confidence judgments have implications for how likely people are to change their mind when encountering the same decision again (Folke, Jacobsen, Fleming, & De Martino, 2016; Kaanders, Folke, & De Martino, in preparation). It is therefore possible that the greater metacognitive accuracy of monolinguals translate into better decisions in the long run. Confidence has also been shown to play a role in collective decision making, where the opinions of more confident people tend to be weighted more heavily (Shea et al., 2014). Because bilingual people tended to report greater confidence on average, this might suggest that bilingual people have relatively greater influence over collectively determined outcomes. Finally, metacognition has been implicated in tasks that require people to choose between staying with a current, known option, versus exploring the environment and trying something new, the so-called exploration-exploitation trade-off (Cohen, McClure, & Angela, 2007; Kolling, Behrens, Mars, & Rushworth, 2012). This implies that the greater confidence in the bilingual population might translates into a greater resilience in sticking with their current

choice, at the cost of being less flexible than their monolingual peers. Obviously, all of these hypotheses are highly speculative and would have to be investigated in future work.

There are a number of limitations of the current study that need to be considered. First, it is conceivable that the lower metacognitive performance by the bilinguals has nothing to do with metacognitive processing but represent a failure of the bilinguals (most of whom had English as a second language) to fully understand the instructions. There is some evidence that is consistent with this account as self-reported second-language language proficiency in the bilinguals correlated with metacognitive ability in Experiment 1 (but not in Experiment 2). However, I find this account unconvincing as the instructions were simple and everyone in both sample had attended or was currently attending university-level education in the UK, which requires a high level of English comprehension. Another, more serious problem with the current study was that there was no direct measure of executive functioning, despite the fact that our original hypothesis was informed by the reported executive functioning advantage in bilinguals. Such a test was not included in the original experiment because of concerns relating to participant fatigue as both the dot discrimination task and the working memory and non-verbal reasoning tasks are quite demanding. It was not included in the second experiment because we wanted to keep the experimental design as similar as possible, while controlling for the potential confound of response times to get as close to a direct replication of the original result as possible. That being said, the lack of a measure of executive functioning is a severe limitation as its inclusion would have allowed to directly explore the relationship between metacognition, executive function and bilingualism, and I would strongly encourage anyone who wants to extend this work to also measure executive functioning directly. It is also too early to tell what these results mean with regards to the cognitive consequences of bilingualism, while we tried to control for obvious confounds such as, age education and non-verbal reasoning ability we did not control for potential socioeconomic differences between the language groups, and second-language ability in the bilinguals were just captured by self-report. Future work in adult samples should control for, or at least measure, a broader range of socio-economic and demographic variables and objectively measure second-language ability so that these factors can be critically examined. The ongoing work on young children that I alluded to earlier in the discussion might provide more insight into how metacognitive ability relates to bilingualism by examining when these differences first manifest, and what other changes might coincide with them. However, while these results do not currently have strong implications for how bilingual cognition differ from monolingual cognition there may be interesting implications for metacognitive research. Specifically, the reported difference in metacognitive performance between the language groups

cannot be explained by differential sensitivity to evidence strength (here operationalised as dot difference) and response time, the two cues that have been most extensively studied in relation to confidence judgements in perceptual decision-making. This suggests that there might be other cues that influence confidence that remain to be discovered, I will explore some such potential cues in the subsequent chapters.

The apparent bilingual disadvantage in metacognition discovered in this work highlight a potential cognitive cost of bilingualism in contrast to the many benefits that have been reported in previous studies. However, this finding seems to provide more question than answers. What is the computational mechanism behind this difference? Is it domain general? When does it manifest? What implications (if any) does it have for higher-order reasoning? All of these are open questions that will hopefully be addressed in future work.

4. Explicit Representations of Confidence Inform Future Value-based Decisions

4.1. Summary

People can report confidence in value-based judgments, but it is not obvious what the benefit of confidence is in the value-domain. In the perceptual domain confidence has been suggested to act as an error signal, which allows people to quickly correct mistakes. In this chapter I apply this idea to the value domain. Two experiments show that confidence judgments predict subsequent changes of mind, when the same options are repeated at a later time. I use a novel computational framework to show that these changes of mind lead to more transitive (internally consistent) choices in participants with high metacognitive ability but not participants with low metacognitive ability. This suggests that confidence is used to improve decision making in value-based judgments as it does in perceptual judgments, and over a longer time-frame than has been tested previously. Previous work has shown that value-based choices are influenced by how long participants look at a given option (dwell time). However, it is unclear if the effect of dwell time depends on the value of the fixated option or not. In both experiments a model that treat dwell-time as an additive boost to evidence accumulation fit the data better than an interactive model, or a null-model where dwell time has no effect. I also identify a novel predictor of confidence judgments from participant eye behaviour leading up to choice, this predictor appears in two snack experiments, but fails to replicate in an experiment with two-armed bandits where the expected value of each arm is learned.

4.2. Introduction

Value-based decisions are ubiquitous in human existence; deciding what to have for lunch, what career to pursue and who to marry are all value-judgments. In order to quantify and study these decisions, economists have long posited a general currency for value that allows us to compare disparate goods (Kahneman & Tversky, 1979), such as buying a candy bar versus renting a movie versus going climbing. With the advent of neuroeconomics there is growing evidence demonstrating that this general value currency is not just a convenient abstraction, but a quantity

represented in the brain, particularly in the ventromedial prefrontal cortex (Levy & Glimcher, 2012).

However, there is another dimension that is almost as central to human decision making as value, but less understood: confidence. Confidence in choices relates to the value of the options in two ways, when options have similar values confidence in the final decision tends to be lower and when one or more options have uncertain values decisions tend to be less confident. The first statement should be intuitively obvious, but the second might require some clarification. Imagine that you are choosing between watching a film and going climbing. You know you enjoy climbing; you have never seen the film under consideration, but a friend you trust recommended it. In the end you decide to go climbing. Despite enjoying yourself (the option you picked had a high value), you are not very confident that you made the best decision because you are uncertain how much you would have enjoyed the film. This example illustrates that confidence captures something about a decision beyond point-estimates of the values of the available options. Recent empirical work shows that confidence only shares a small amount of variance with value (De Martino et al., 2013), and that self-reported confidence predicts choice behaviour above and beyond self-reported value of the options.

Confidence is an explicit measure of uncertainty in value-based choice, but there are also other, implicit measures that imply uncertainty in the choice process. The strength of the evidence favouring a specific option and response time are both associated with the level of uncertainty in the choice process and tends to be associated with confidence judgments (Baranski & Petrusic, 1994, 1998; De Martino et al., 2013; Festinger, 1943; Kiani & Shadlen, 2009; Ratcliff & Starns, 2009; Vickers & Packer, 1982). Kiani, Corthell and Shadlen (2014) recently showed that the relationship between response time and confidence is not merely incidental but that response times causally influence confidence (Kiani et al., 2014).

Eye tracking provides a promising avenue for finding additional behavioural markers of uncertainty, as fixations have been shown to influence the accumulation of evidence in value-based choices. Specifically, Rangel and colleagues have suggested that fixating on an item allows the value of that item to be sampled more efficiently (Krajibich & Rangel, 2011; Krajibich, Armel, & Rangel, 2010). This account has since been challenged by Cavanagh and colleagues who suggest that looking at an object boosts the evidence accumulation in favour of that object, independent of its value (the evidence accumulation is additive rather than interactive; Cavanagh, Wiecki, Kochar, & Frank, 2014). To the best of my knowledge Cavanagh and colleagues are the only ones who have directly compared these models, but they examined this effect on stimuli

with learned values. Here, I will directly compare these accounts in snack data, which is closer to Krajbich and Rangel's original research designs. I will also use eye tracking data to explore the relationship between eye movements and explicit confidence, something that, to the best of my knowledge, has not been attempted before this study.

The computational underpinnings of confidence are still debated. One popular suggestion is that confidence captures noise in the stochastic evidence accumulation process leading up to a decision (De Martino et al., 2013; Kepecs et al., 2008; Kiani et al., 2014). Lebreton and colleagues have suggested that confidence is an inherent property of value computations, and that confidence has a quadratic relationship to the value signal (Lebreton, Abitbol, Daunizeau, & Pessiglione, 2015, see also Barron, Garvert, & Behrens, 2015). Several studies suggest that confidence for both value-based and perceptual decisions might be represented in the rostro-lateral prefrontal cortex (De Martino et al., 2013; Fleming, Huijgen, & Dolan, 2012; Fleming et al., 2010; Rounis et al., 2010).

But why spend neural resources representing confidence? How does an explicit representation of confidence benefit the agent? One suggestion is that confidence allows people with different beliefs to quantify the relative strength of their beliefs, which might improve group decision making (Bahrami et al., 2012; Bang et al., 2014; Shea et al., 2014). Others have proposed that confidence judgments might be useful to the individual by helping to guide future decisions (Lau & Rosenthal, 2011). Proponents of this view suggest that confidence seems to improve learning under uncertainty (Meyniel, Schlunegger, & Dehaene, 2015), that confidence judgments and error correction seem to be driven by the same computations (Resulaj et al., 2009; Yeung & Summerfield, 2014), and that changes of mind and confidence judgments seem to operate on the same evidence scale (Van Den Berg et al., 2016).

An open question regarding the utility of confidence is whether confidence judgments in a current choice can influence future decisions. Returning to the example from the beginning of the introduction: given that you weren't very confident that climbing was the most enjoyable activity, are you more likely to pick the film when a similar choice comes up days later?

Additional research from our lab indicates that you might. In another study, we allowed people to reverse a perceptual decision following a confidence judgment and to sample additional information about the chosen and unchosen options. We found that low-confidence judgments tended to predict changes of mind, and that people who are more metacognitive use their confidence judgments to sample the available information more effectively, which in turn leads to better decisions (Kaanders, Folke, & De Martino, in preparation).

The work presented here extends these findings by allowing time to pass between the initial choice and the potential change of mind. Here, people first chose between a set of snack items and subsequently gave a confidence judgment. Sometime later during the experiment, they would see the same options again (with positions counterbalanced) and be asked to repeat the choice. By predicting changes of mind from confidence in this context, I move away from an immediate error correction framework and suggest that confidence in a choice might have implications for long-term planning.

Another key difference between the perceptual research cited above and this work is that perceptual decisions have correct and incorrect options, so one can legitimately ask if confidence judgments help improve decision making or not. There is no obvious measure of accuracy in the value domain because value judgments are, by their very nature, subjective. One way to approximate correctness is internal consistency: if I prefer an apple over a banana and a banana over a pear I should prefer an apple over a pear, or more formally and generally if $A > B$ and $B > C$ then $A > C$. This form of internal consistency is known as transitivity and is a normative prescription in utility theory (Von Neumann & Morgenstern, 2007). The reason that transitive preferences are preferable is that they are rational; intransitive preferences lead to situations where a set of individually agreeable trades lead to a net loss. Say that an agent prefers apples to oranges, oranges to pears and pears to apples. The agent currently has an orange. It is now possible for a trader to offer an apple in exchange for the orange and some arbitrarily small amount of money (say one pence), then offer to trade the apple for a pear and another pence and finally offer to trade the pear for an orange and a pence. Given their preference structure the agent should accept all three trades and end up three pence poorer with the same orange they started with. What is worse, this transaction cycle could be repeated an arbitrary number of times creating arbitrarily large losses (See the dutch-book problem; Hájek, 2008). However, failure of transitivity has been reported in human choices and represents an exemplar violation of economic rationality and, more generally, of logical consistency (Camerer & Ho, 1994; Loomes, Starmer, & Sugden, 1991). In this research I used a novel method to generate the ordered ranking of items corresponding to the most transitive choices for each participant. I then tagged the choices that violated these optimal rankings (transitivity violations), which allowed me to explore how confidence relates to internal consistency.

This chapter will provide a number of novel insights. First it will examine what factors contribute to the construction of confidence during the formation of a value-judgment, being the first work to include eye behaviour in such models. Second, it will directly compare additive and

interactive models of the role of fixations in value-based decision making. Third, it will explore how explicit representations of confidence subsequently inform and improve future decisions.

4.3. Methods

4.3.1. Experimental Procedures, Experiment 3

Participants were required to make binary choices between 16 common snack items. Participants were asked to choose between each combination of the items ($n = 120$) twice, counterbalanced across the left-right spatial configurations (total number of choices = 240). After each choice, participants indicated their confidence in their decision on a continuous rating scale. The edges of the confidence scale was labelled “low” and “high” and participants were asked to give their confidence that they had picked their preferred option. No feedback was given on how participants used the confidence scale. Neither choices nor confidence ratings were time constrained. The second presentation of the same pair was randomly interleaved during the task with the only constraint being that the same pair was never repeated in immediately subsequent trials. Participant eye movements were recorded throughout this task. At the end of the experiment, one choice from this phase was played out and the participant had the opportunity to buy the chosen item by means of an auction administered according to the Becker-DeGroot-Marschak (BDM) procedure (Becker, DeGroot, & Marschak, 1964). This procedure encouraged participants to choose preferred snacks during the eye tracking choice phase since they only had a chance to win snacks they chose. Once a snack had been selected the experimenter randomly extracted a price from a uniform distribution (£0 to £3)—the ‘market price’ of that item. If the participant’s bidding price (willingness-to-pay) was below the market price, no transaction occurred. If the participant’s bidding price was above the market price, the participant bought the snack item at the market price. At the end of the experiment, participants had to remain in the lab for an additional hour. During this hour, the only food they were allowed to eat was the item purchased in the auction, if any. At the end of the waiting period participants were debriefed and thanked for their participation. Participants were paid £25 for their time, deducting the cost of the food item, if they bought any. Both tasks were programmed using MATLAB 8.0 (MathWorks) running the Psychophysics toolbox (<http://psychtoolbox.org>) as well as the Eyelink toolbox extensions (Brainard, 1997; Cornelissen, Peters, & Palmer, 2002).

4.3.2. Experimental Procedures, Experiment 4

Participants gave their willingness to pay for 72 common snack food items on a scale ranging from £0-£3, in a BDM procedure, similar to the one in Experiment 1. Next they completed a

choice task where, in each trial, they had to pick their favourite item out of three options. The triplets presented in the choice task were tailored for each participant from their willingness-to-pay ratings. The items were divided into high-value and low-value sets by a median split. The 36 high-value items were randomly combined into 12 high-value triplets; this procedure was mirrored to generate 12 low-value triplets. The high-value items and low-value items were then mixed to generate medium value triplets, so that there were 12 triplets consisting of two high-value items and one low-value item, and 12 triplets showing the reverse ratio. This resulted in 48 unique triplets, with counterbalanced spatial configurations (total trials = 144), split into three blocks. The trial order was pseudo-randomised with the constraint that the same triplet was never shown twice in a row. In the subsequent choice task, the triplets were presented inside 3 white squares in an equidistant 2x2 grid (one position on the grid was left empty, this was randomly determined). I used a gaze-contingent paradigm in which the items were only visible when the participant made a fixation inside one of the squares, so that the participant could only see one item at a time. Participants had unlimited time to make up their mind and could make as many fixations as they wished. After each choice, participants indicated their confidence in their decision on a visual analogue rating scale, similar to the one used in Experiment 3 (without any time constraints). Participants' eye movements were recorded throughout the choice task. Both the choice task and the willingness to pay procedure were programmed in Experiment Builder version 1.10.1640, SR-Research, Ontario. Following the choice task, an auction based on the BDM-ratings took place (see Experiment 1). After the auction, participants had to remain in the lab for an additional hour, as in Experiment 1. At the end of the waiting period participants were debriefed and thanked for their participation. Participants were paid £15 for their time, deducting the cost of the food item, if they bought any.

4.3.3. Experimental Procedures, Experiment 5

Participants completed a two-armed bandit task with the aim to maximise their earnings. Each bandit had the same reward magnitude (2.08 pence) but the rate of reward differed. In each trial two bandits were shown and participants selected one with a key press. Following the choice a circle appeared at the centre of the screen to cue participants to fixate on the centre prior to receiving feedback. Once participants fixated on the central circle it disappeared and boxes appeared around the bandits. The boxes were either orange or purple, and the colour indicated whether the bandit was rewarded or not during that trial (the colour associated with reward was counterbalanced between participants). The reward probabilities of the bandits were independent, so for any given trial both bandits could win, or no bandit could win, or one could

win but not the other. Participants' eye movements were recorded both during the choice phase and the feedback phase of each trial. For a visual representation of the task structure see Figure 4.1. Participants first completed 20 practice trials with two bandits with stable reward rates (0.75 and 0.25 respectively) in order to familiarise them with the task structure and to make sure they understood the aim of the task. The practice trials did not count towards the participant earnings.

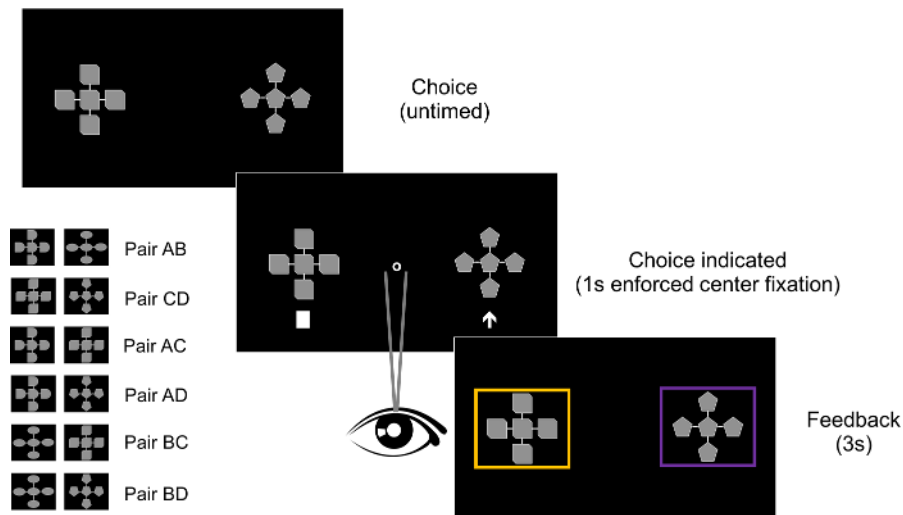


Figure 4.1. Task Structure for Experiment 5

Participants were asked to choose between one of two symbols on each trial, each with their own independent probability of yielding a reward. Choices were untimed and indicated with a button press. After the choice was made, feedback was displayed for each item after the participant fixated on a center marker for one second. The color of the box indicated the presence or absence of reward for both options (note that the reward rate of each item was independent). The color assigned to each outcome was counterbalanced across participants. Four different symbols were arranged in six unique pairs, and each pair was presented 40 times for a total of 240 trials. Left-right presentation of each pair was counterbalanced within blocks of 12 trials.

Then followed 240 main trials. These trials had four different bandits, making six unique bandit pairs that were shown 40 trials each. The screen positions of the bandits were counterbalanced for each pair. The four bandits started with reward rates of 0.2, 0.4, 0.6 and 0.8 respectively. The reward rates were randomly assigned to the bandits at the beginning of each experimental session. The reward rates for the main bandits were not stable throughout the experiment. Each time a bandit was shown its reward rate changed by 0.03. If the reward rate of a bandit reached 0 or 1, it “bounced” to 0.03 or 0.97 but other than that it drifted randomly. Participants were

aware that the reward rates of the bandits would change over time. The bandit task was programmed in Experiment Builder version 1.10.1640, SR-Research, Ontario.

4.3.4. Participants, Experiment 3

29 participants took part in the study. One participant was excluded because the BDM estimates were poor predictors of their choice (choose the high valued item 63% of the time, compared to the group mean of 83%, $sd=7\%$). Thus 28 participants were included in the analysis (13 females, age: 19-73). All participants were required to fast four hours prior to taking part in the experiment. Blood glucose levels were taken to test their adherence to this criterion (mean glucose level = 83.57mg/dl, $sd = 10.90\text{mg/dl}$).

4.3.5. Participants, Experiment 4

30 participants completed the study. Of these thirty, three were excluded because of a limited range in their BDM ratings (they gave the exact same price for more than 28% of the items). An additional three participants were excluded for limited range in using the confidence scale (reporting the same level of confidence in more than half of the trials). 24 participants were included in the main analyses (17 females, age: 21-38). All participants were required to fast for four hours prior to doing the experiment.

4.3.6. Participants, Experiment 5

30 people participated in this study. One participant was excluded because their choices did not reflect the reward rates of the bandits for the second half of the main trials (a logistic model predicting choice by difference in reward rates did not predict choice better than an intercept-only model, $p = .70$). Thus the final sample included 29 people (20 female) with a mean age of 26 ($sd = 6.37$). All participants gave informed consent and were paid a £10 show up fee and up to an additional £5 based on their performance.

4.3.7. Eye Trackers

For Experiment 3, eye gaze was sampled at 250 Hz with a head-mounted SR Research Eyelink II eye tracker (SR-Research, Ontario). For Experiments 4 and 5, eye movements were recorded at 1000Hz with an EyeLink 1000 Plus (SR-Research, Ontario).

4.3.8. Preparation of the Eye Data, Experiment 3

Areas of Interest (AI) were defined by splitting the screen in half, creating two equal sized areas. Fixation on the left AI was assumed to be directed towards the left snack item, and vice versa. I

constructed two variables from the eye tracking data: the difference in dwell time between the two AIs (DDT), and gaze shift frequency (GSF). DDT was calculated by subtracting the total dwell time on the left side from the total dwell time on the right side. GSF was a count of how many times participants shifted their gaze from one AI to the other during a given trial.

4.3.9. Preparation of the Eye Data, Experiment 4

AIs were pre-defined by the 3 squares that participants had to fixate to view the items. I derived four variables from the eye tracking data: the total dwell time in each of the three AIs for a given trial, and GSF. Following Experiment 3, GSF measured the number of fixations in one AI immediately followed by a fixation in another AI. To ensure that participants paid attention, I excluded trials where participants did not fixate on every option available at least once. 13 trials out of 3457 were excluded from the analysis for this reason.

4.3.10. Preparation of the Eye Data, Experiment 5

AIs were predefined by the 2 squares on the screen (400 x 329 pixels) covering the bandits. I derived three variables from the eye tracking data: the total dwell time in each AI for a given trial, and GSF. In line with the two snack experiments, GSF measured the number of fixations in one AI immediately followed by a fixation in another AI.

4.3.11. Hierarchical Models

Note that all predictors entered into the hierarchical models are z-scored on the participant level, and that response time was log-transformed prior to being z-scored to make RT distributions approximately normal. All models reported in this chapter allowed for random intercepts and random slopes at the participant level. The individual difference analyses investigating change of mind and transitivity did not depend on hierarchical parameter estimates but unpooled estimates. The rationale behind this choice was that for both analyses I was interested in studying between-participant variations (Fig. 4c and Fig 5c) that could be potentially affected due to shrinkage of parameters towards the group mean that characterise hierarchical models (Gelman & Hill, 2006).

4.3.12. Drift Diffusion Models

I fitted three hierarchical drift diffusion models (DDM) to the choice data from Experiment 3 to evaluate whether fixation time influenced evidence accumulation in an additive or interactive fashion (see Chapter 2 for details on the software and fitting procedure). Recall that DDM's have four important free parameters: drift rate (v), boundary separation (a), the starting point (z) and

the non-decision time (t). In order to make the models as similar to the models in the Cavanagh paper as possible I predicted the probability of picking the option with the highest BDM-value (trials where both options had the same value, about 3% of the dataset were excluded from these analyses). Boundary separation and non-decision time was allowed to vary between participants, whereas the starting point was set to 0.5 for everyone. The three models differed in how they calculated drift rates:

$$Model_0: v = \beta_0 + \beta_1 \times (BDM_{high} - BDM_{low})$$

$$Model_a: v = \beta_0 + \beta_1 \times (BDM_{high} - BDM_{low}) + \beta_2 \times (fixation\ ratio_{high} - fixation\ ratio_{low})$$

$$Model_i: v = \beta_0 + \beta_1 \times (fixation\ ratio_{high} \times BDM_{high} - fixation\ ratio_{low} \times -BDM_{low}) + \beta_2 \times (fixation\ ratio_{low} \times BDM_{high} - fixation\ ratio_{high} \times -BDM_{low})$$

Where β_0 was a random intercept parameter that was allowed to vary by participant and β_1 and β_2 were sensitivity parameters that were fixed across the entire sample. BDM_{high} was the BDM value in pounds of the highest-valued option of that trial and BDM_{low} was the lowest valued option. Fixation ratio_{high} was the proportion of the response time spent looking at the highest valued option and fixation ratio_{low} was the proportion of response time spent looking at the low-value option. Note that the interactive model is a simplification of Krajbich and Rangel's original aDDM because the original modelled the evolution of the drift throughout the trial and therefore accounted for the order of the fixations, here I model the drift rate, the average slope of the path of the particle throughout the trial, and ignore the order of the fixations, a simplification first introduced by Cavanagh and colleagues (2014).

4.4. Results

4.4.1. Relation Between Confidence and Choice

First I tested whether confidence judgments related to accuracy in the value domain, by testing if high-confidence choices corresponded more with the participants stated preferences than low-confidence choices. This was evaluated with a hierarchical logistic regression model that predicted the probability of the participant choosing a reference item (the right item in Experiment 3, the first item encountered based on western reading order in Experiment 4). Choices were predicted by difference in value (DV), operationalised as the difference in BDM value between the reference item and the second item (Experiment 3) or the mean of the other

two items (Experiment 4), confidence, summed value (SV; the total BDM value of all the options in that trial), difference in dwell time (DDT) and the interactions between DV and Confidence and DV and SV. In line with a wealth of previous research (Boorman, Behrens, Woolrich, & Rushworth, 2009; FitzGerald, Seymour, & Dolan, 2009; Lebreton, Jorge, Michel, Thirion, & Pessiglione, 2009; Levy, Lazzaro, Rutledge, & Glimcher, 2011; Plassmann, O'Doherty, & Rangel, 2010) I found that difference in value (DV) was a reliable predictor of participants' choices in both of the snack experiments (hierarchical logistic regression; Experiment 3: $z = 11.48$, $p < 10^{-10}$, Figure 4.2. c & f; Experiment 4: $z = 6.66$, $p < 10^{-9}$, Figure 4.2. b & e). Both studies also showed a significant negative interaction between the summed value of all options (SV) and value difference (DV) (Experiment 3: $z = -3.08$, $p = .002$; Experiment 4: $z = -2.84$, $p = .005$), indicating that participants were more able to use DV to guide their choice when item values were low, compared to when items were high in value (Figure 4.2. c & f). To my knowledge this effect has never been reported before but is consistent with the Weber–Fechner law in sensory perception, in which the resolution of percepts diminishes for stimuli of greater magnitude and it is compatible with the notion of divisive normalisation (Carandini & Heeger, 2012; Louie, Khaw, & Glimcher, 2013; Soltani, De Martino, & Camerer, 2012). Confidence, unlike DV, was not in itself a predictor of choice (right or left item) but instead correlates with choice accuracy. Both experiments replicated this finding, which had previously been reported in De Martino et al. (2013). This effect is shown here (Figure 4.2. b & e) using a logistic regression model predicting choice. I found a significant positive interaction between DV and confidence (Experiment 3: $z = 7.38$, $p < 10^{-10}$; Experiment 4: $z = 5.78$, $p < 10^{-8}$), with DV predicting choice more strongly for trials when confidence is high.

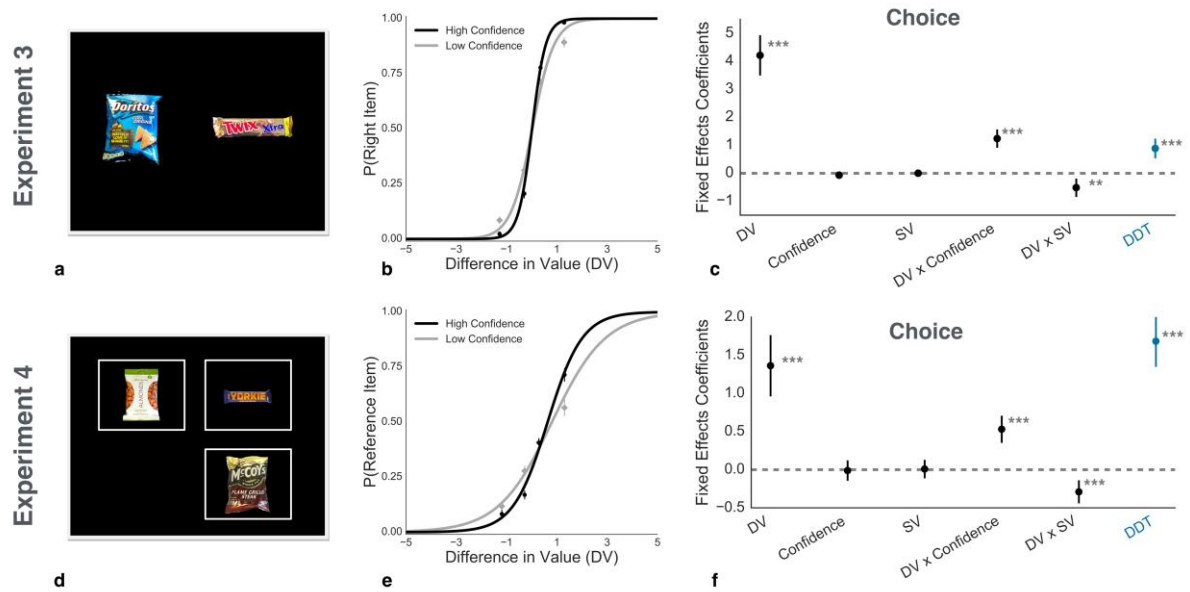


Figure 4.2. Relation Between Confidence and Choice

Eye tracking tasks: **(a)** In Experiment 3, participants were presented with two snack items and were then required to choose one item to consume at the end of the experiment. **(d)** In Experiment 4, participants chose between three options, and the presentation of the stimuli was contingent on which box participants looked at. In both experiments, participants indicated their confidence that they had made a correct decision on a visual analogue scale, after each choice had been made. **(b)** Probability of choosing the item on the right as a function of the difference in value between the options. **(e)** Probability of choosing the reference item (see paragraph 4.4.1.), as a function of the standardised value difference between the reference item and the mean value of the alternatives. Black line: high confidence trials, grey line: low confidence trials (as determined by a median split). **(c and f)** Fixed effects coefficients from hierarchical logistic regression models predicting choice (DV= difference in value; SV= summed value; DDT= difference in dwell time, DV x Confidence= Interaction of difference in value and confidence; DV x SV= Interaction of difference in value and summed value). The graph for Experiment 3 (C) shows the coefficients predicting the probability of choosing the right-hand option; the graph for Experiment 4 (F) shows the coefficients predicting the probability of choosing the reference option. Error bars show 95% CIs. *** = $p < .001$; ** = $p < .01$; * = $p < .05$.

4.4.2. Dynamics of Information Sampling

I then examined the dynamics of eye movements between items during the choice – both the total amount of time participants spent looking at each item and also how frequently gaze shifted

between items. Replicating previous studies (Krajbich & Rangel, 2011; Krajbich et al., 2010) I found that gaze behaviour correlated with choice. In line with Cavanagh and colleagues I found that difference in fixation time and difference in value predicted choice additively (Cavanagh et al., 2014). I compared the Krajbich and Cavanagh accounts of the influence of fixation time by fitting three drift diffusion models to Experiment 3, a null model that predicted drift rate only from the value difference between the options, an additive model that included an additive boost to drift rate for the option that had been fixated on longer, and an interactive model that weighted the value of each option with how long it had been fixated (see the methods for details). I found that the additive model explained the data better than the null model and the interactive model ($DIC_0=19\,704$, $DIC_a=19\,305$, $DIC_i=19\,989$). Surprisingly, the null model fit the data better than the interactive model. For the winning additive model both difference in value and difference in fixation ratio positively predicted drift rate ($\text{coefficient}_{\text{value}}=0.80(0.02)$, $\text{coefficient}_{\text{dwellratio}}=0.79(0.04)$).

I also constructed a novel measure that captured the dynamics by which information was sampled. This new measure, called ‘gaze shift frequency’ (GSF), indexed how frequently gaze shifted between the options presented on the screen. This measure is independent of difference in dwell-time: for a constant allocation of time between the options (e.g. 3 seconds for the left-hand option and 5 seconds for the right-hand option) one may shift fixation only once (e.g. switching from left to right after 3 seconds have elapsed; low gaze shift frequency) or shift many times between the two options (high gaze shift frequency). In practice DDT and GSF appears to be weakly correlated, both in Experiment 3 ($r=.10$, $p<10^{-10}$), and in Experiment 4 ($r=.06$, $p<10^{-4}$). GSF was positively correlated with response time in Experiment 3 ($r=.57$, $p<10^{-10}$) and Experiment 4 ($r=.63$, $p<10^{-10}$), suggesting that participant moved their gaze more during longer trials. For the learning experiment, I found no correlation between GSF and DDT ($r=-.003$, $p=.79$), and only a weak relationship between GSF and RT ($r=.05$, $p<10^{-4}$).

In order to test how these two gaze metrics influenced behaviour I ran a set of hierarchical logistic regression models that predicted choices from DV, DDT and GSF and interactions between DV and DDT and DV and GSF for Experiments 3, 4, and 5 (see Figure 4.3.). In Experiment 3 I found a double dissociation between the impact of difference in dwell time and gaze shift frequency on choice. In other words, dwell time produced a bias in choice (shift in the psychometric function; $z=16.95$, $p<10^{-10}$) but did not affect choice accuracy (slope in the psychometric function; $z=-0.87$, $p=.38$). On the contrary, gaze shift frequency had no effect on choice *per se* ($z=1.74$, $p=.08$) but negatively predicted the accuracy of the choice ($z=-4.17$, $p<10^{-4}$). The results from Experiment 4 were noisier. As in Experiment 3, GSF negatively predicted

accuracy ($z=-3.71$, $p<10^{-4}$), but it also produced a bias in choice ($z=-2.46$, $p=.01$), this bias in choice is probably an artefact of how GSF interacted with DDT, because if DDT was excluded from the model GSF did not predict choice on its own (but it still predicted choice in interaction with value, see Appendix 1). DDT had a strong main effect on choice ($z=21.19$, $p<10^{-10}$), but also interacted with value difference ($z=-2.29$, $p=.02$). However, while the dissociation was not complete the main effect was much stronger for DDT than GSF, and the interaction effect was much stronger for GSF than for DDT. I reanalysed the data by simply predicting the probability of picking the highest valued item; here GSF predicted choices above unsigned value difference ($z=-6.24$, $p<10^{-9}$), but unsigned DDT did not ($z=-0.65$, $p=.52$), in line with the idea that GSF correlates with choice accuracy but DDT does not. For the learning experiment, DDT influenced only bias (as in Experiment 3; $z=16.92$, $p<10^{-10}$) and GSF did not predict bias ($z=-1.16$, $p=.24$) or accuracy, $z=0.34$, $p=.73$).

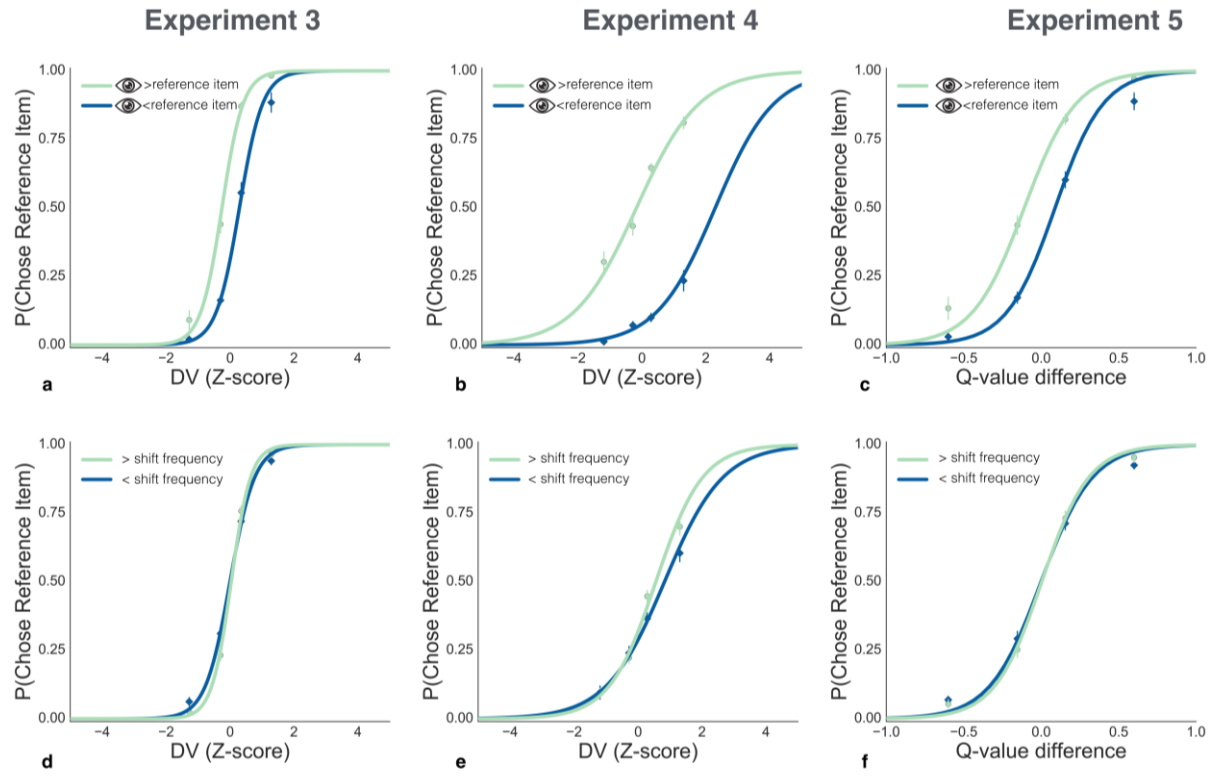


Figure 4.3. Dynamics of Information Sampling

The plots show data from Experiments 3, 4, and 5 illustrating the differential effects of each eye tracking parameter on choice. The y-axes show the probability of choosing the reference item and the x-axes show difference in value (z-scored) for Experiments 3 and 4 and differences in Q-values for Experiment 5. Data are split into values above the median (green line) and below the median (blue line) for difference in dwell time (**a-c**) and gaze shift frequency (**d-f**). Points

represent quartiles of DV. Error bars show standard errors. Difference in dwell time is associated with a shift of the logistic fit with no change in slope, consistent with a biasing effect on choice (i.e. people are more likely to pick the item they look at for longer, independent of difference in value). In contrast, gaze shift frequency modulates choice accuracy (change in slope) but is not associated with a bias in choice in the decision-experiments. Gaze shift frequency does not appear to influence choice in any way during the learning experiment.

4.4.3. Factors that Contribute to Confidence

I then investigated which variables contributed to the subjective representation of confidence during value-based choice. Previous work has shown an interrelationship between absolute difference in value ($|DV|$), response time (RT) and confidence (i.e. participants are more confident both when $|DV|$ is high and their choices are faster; De Martino et al., 2013). These findings are in line with the conceptual relation between confidence, difference in stimulus strength (indexed by $|DV|$ in the value-based framework) and RT (Kiani et al., 2014; Kiani & Shadlen, 2009). Using hierarchical linear regression models without intercepts, I observed this same relation in the current study. In both experiments I found that $|DV|$ was a significant predictor of confidence (Experiment 3: $t=12.62$ $p<.10^{-10}$; Experiment 4: $t=8.12$, $p<10^{-7}$). I also found that RT was a negative predictor of confidence (Experiment 3: $t=-11.00$, $p<10^{-10}$; Experiment 4: $t=-7.57$, $p<10^{-6}$). Additionally, I found that summed value positively predicted confidence, meaning that participants tended to be more confident when the options were all high in value Experiment 3: $t=3.58$, $p=.001$; Experiment 4: $t=4.77$, $p<10^{-4}$). This finding indicates that overall value might boost confidence.

I also included the eye variables GSF and $|DDT|$ in the model, to investigate how they related to confidence judgments. While $|DDT|$ was a weak positive predictor of confidence in Experiment 3 ($t=2.15$, $p<.05$) it did not predict confidence in Experiment 4 ($t=-0.52$, $p=.60$). GSF was a robust negative predictor of confidence in both Experiments 3 and 4 (Experiment 3: $t=-3.44$, $p=.002$; Experiment 4: $t=-7.41$, $p<10^{-6}$) see Figure 4.4. a and b. In other words, in trials in which participants shifted their gaze more often between the available options their confidence was lower, even accounting for changes in $|DV|$ SV and RT. The four-way relationship between $|DV|$, RT, GSF and confidence is plotted in Figure 4.4. c and d.

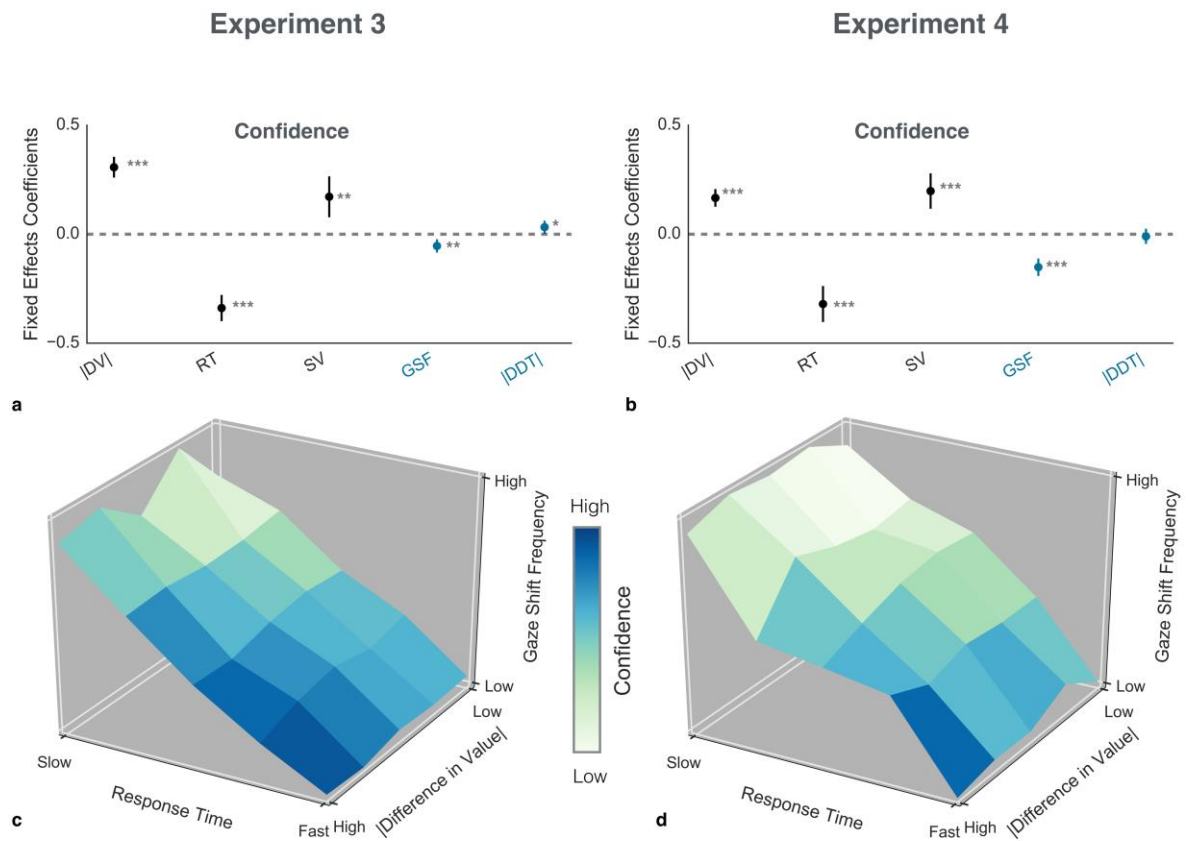


Figure 4.4. Factors that Contribute to Confidence

(a-b) Graphs show coefficient plots for the fixed-effect coefficients in hierarchical regression models predicting confidence for Experiments 1 and 2, respectively. Error bars show 95% CIs. *** = $p < .001$; ** = $p < .01$; * = $p < .05$. ($|DV|$ = absolute difference in value; RT = reaction time; SV = summed value; GSF = Gaze Shift Frequency, $|DDT|$ = absolute difference in dwell time). **(c – d)** 4-D heat maps showing mean z-scored confidence as a function of subject specific quantiles of response time, absolute difference in value and gaze shift frequency. Gaze shift frequency and confidence are both influenced by response time and the absolute value difference between the options.

4.4.4. Confidence Predicts Change of Mind

In the two snack experiments participants saw the same exact choice sets on more than one occasion. In Experiment 3 each pair was presented twice; in Experiment 4 each triad was presented three times (counterbalancing for different spatial locations). This design allowed me to determine factors affecting changes of mind when the same options are encountered again in a subsequent trial. Note that the way I define change of mind above is different from how it is often defined in perceptual decision making, namely as a choice reversal within the same trial

after further processing of sensory information (Bronfman et al., 2015; Moran et al., 2015; Resulaj et al., 2009; Van Den Berg et al., 2016). The hypothesis I sought to test was that an explicit representation of uncertainty in a choice (reported as confidence) would influence behaviour when the same options were presented again during a different trial. In a hierarchical logistic regression, lower confidence at the first presentation was indeed associated with increase in change of mind at the following presentation, in both Experiments 3 and 4 (Experiment 3: $z=-6.16$, $p<10^{-9}$; Experiment 4: $z=-5.21$, $p<10^{-6}$). The effect of confidence in predicting change of mind remained robust after controlling for all the other factors that might correlate with the stability of a choice such as $|DV|$ and RT. Notably none of the eye tracking measures played a significant role as predictor of change of mind when included in the regression analysis (Fig. 4 *coefficients in blue*). Note that this was still true when confidence was excluded from the regression analysis ($GSF_{\text{Experiment 3}}$: $z=-0.68$, $p=.49$; $|DDT|_{\text{Experiment 3}}$: $z=-0.32$, $p=.75$; $GSF_{\text{Experiment 4}}$: $z=0.59$, $p=0.55$; $|DDT|_{\text{Experiment 4}}$: $z=-0.13$, $p=0.90$).

This is particularly interesting for GSF because of its significant negative relation with confidence (see Fig.3). This result suggests the hypothesis that the low level (and possibly implicit) measure of uncertainty gathered by GSF is insufficient to trigger a delayed change of mind. On the contrary, an explicit representation of uncertainty (expressed through confidence) allows individuals to capitalise on their ‘knowledge about their ignorance’ and make a difference choice when similar options are presented later.

It is important to note that just because confidence predicted changes of mind in subsequent trials, it does not necessarily follow that low confidence judgments are causing subsequent changes of mind. We know that low confidence judgments are associated with a noisier decision-process, leading to more decisions that violating the preferences of the participant. So perhaps confidence simply tags noisy decisions as part of some error monitoring process (Yeung & Summerfield, 2012, 2014). When the same choice repeats the decision process is less noisy because of recursion to the mean and the less noisy decision process causes the highest-value option to be chosen. This results in a change of mind in relation to the first low-confidence choice, but confidence is not having any causal influence. I attempted to account for this by adding a dummy-variable coding whether the highest-value option was chosen in the original trial, and confidence still predicted future changes of mind in Experiment 3 ($z=-4.79$, $p<10^{-5}$) and Experiment 4 ($z=-4.88$, $p<10^{-5}$). Therefore, it seems probable that confidence causally influence future changes of mind.

Next, I examined whether individual differences in metacognition related to changes of mind. I reasoned that the impact of confidence on changes of mind would be more prominent in participants who have enhanced metacognitive skills, i.e. those whose explicit confidence ratings more accurately track the level of uncertainty underlying their decision process. In order to test this hypothesis I calculated an individual index of metacognitive sensitivity by computing the difference in slope between psychometric functions on high and low confidence trials (De Gardelle & Mamassian, 2014; De Martino et al., 2013; Fleming & Lau, 2014). I then ran a logistic model to predict changes of mind at later presentations using confidence measured at earlier presentations with the same stimuli. In line with my initial hypothesis, the impact of confidence on changes of mind is stronger in those participants with greater metacognitive accuracy ($r = -0.35$, $p = 0.01$; Figure 4.5.). Note that the relationship is negative because the influence of confidence on changes of mind should be negative (so that changes of mind are more probable when confidence in the initial choice is low).

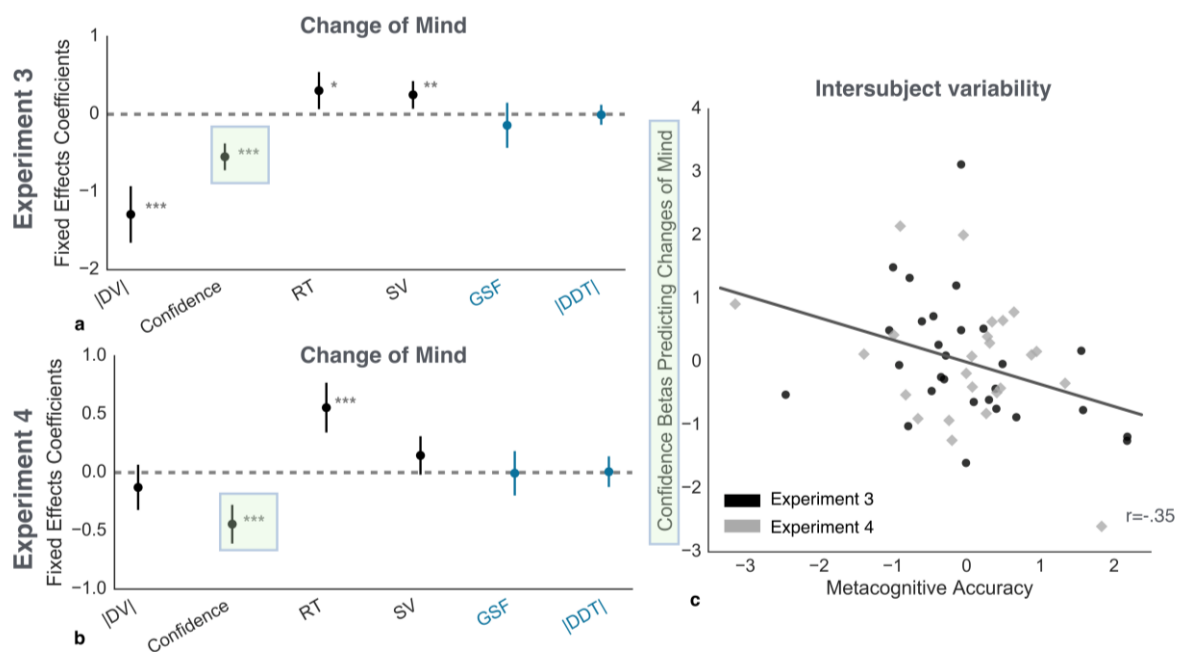


Figure 4.5. Confidence Predicts Change of Mind

(a-b) Coefficient plots for the fixed effects coefficients from hierarchical logistic regression models predicting future changes of mind. Error bars show 95% CIs. *** = $p < .001$; ** = $p < .01$; * = $p < .05$. (|DV| = absolute difference in value; RT = reaction time; SV = summed value; GSF = gaze shift frequency; |DDT| = absolute difference in dwell time) (c) Correlation between metacognitive accuracy and the coefficients for confidence ratings predicting future changes of mind (highlighted in pale green). Participants with greater metacognitive accuracy are

more likely to change their mind following a low-confidence judgment; note that the correlation is negative because the relationship between confidence and changes of mind is itself negative (lower confidence increases the probability of subsequent changes of mind). Participants from Experiment 3 are represented by black dots, participants from Experiment 4 are represented by grey squares. Both axes (x and y) are z-scored for each experiment separately.

4.4.5. Link Between Confidence and Transitivity

In the analysis presented in the previous paragraph I established a link between confidence judgments and changes of mind. A change of mind is agnostic to the quality of the decision because the result might be better or worse than the original choice. However, as mentioned in the introduction, not all choices are equal; some of them can result in a more closely consistent pattern of choices than others, a parameter that can be indexed measuring the level of transitivity across all decisions. In order to test the relation between confidence and transitivity I estimated the (idiosyncratic) preference ranking of items that led to the lowest number of transitivity violations. Note that it is extremely complex to find an optimal ranking order for choice sets with more than a handful of items; however, a number of efficient algorithms that approximate a numerical solution have been developed for pairwise comparisons. Here I used the Minimum Violations Ranking (MVR) algorithm (Pedings, Langville, & Yamamoto, 2012) that minimizes the number of inconsistencies in the ranking of the items shown to each participant. This method enabled me to tag choices as transitivity violations (TV) of the optimal ranking calculated via the MVR algorithm. Because most of these methods are not suited for ternary choice the analyses presented in this section were performed only on data collected for the experiment using binary choice (Experiment 3). After having ordered the participants' choices according to the MVR algorithm, 4.5% of all decisions were classified as transitivity violations. I then split the dataset into trials in which participants reported high confidence and trials in which they reported low confidence (median split). Most of the transitivity violations took place during low confidence trials (85% of the transitivity violations) as opposed to high confidence trials (15% of the transitivity violations). While these results are consistent with previous evidence provided in this paper and elsewhere (De Martino et al., 2013), it is important to highlight that for this analysis did not rely on BDM value estimates and that my approach to generate the optimal ranking did only include choice trials. These show that that the link between confidence and the quality of a value-based decision is robust independently of the method used to assess quality. In order to test which factors accounted for transitivity violations on a trial-by-trial basis I constructed a set of hierarchical logistic regression models. Absolute difference in value

($|DV|$) was a robust negative predictor of TV ($z=-6.41$, $p<10^{-9}$) meaning that participants were more likely to violate transitivity during trials when the items were closer in value. Critically, this same model showed that even when $|DV|$ was accounted for, confidence was a negative predictor of transitivity violations ($z=-6.25$, $p<10^{-9}$). In other words, participants felt less confident during those trials in which they went against their best-fitting preference order. Both response time ($z=4.17$, $p<10^{-4}$) and summed value ($z=2.46$, $p=.01$) positively predicted transitivity violations: trials in which the value of both options was higher and/or in trials in which their responses were slower, participants' choices were more likely to result in transitivity violations. Note that this is another metric showing that summed value appears to be associated with lower quality choices, despite being a positive correlate of confidence. Finally, GSF negatively predicted transitivity violations value ($z=-2.56$, $p=.01$), meaning that participants shifted their gaze less during trials that violated transitivity when the other variables were accounted for. This negative relationship between GSF and transitivity violations is surprising, because I suspect that GSF is a low level measure of choice uncertainty. However, because of the high-correlation between RT and GSF this coefficient estimate might not be reliable (Gelman & Hill, 2006). Indeed, removing RT from the model causes GSF to become insignificant ($z=0.21$, $p=.83$), but does not influence the other predictors. Additionally, decoupling GSF from RT by making it into a rate (GSF/sec) also makes it an insignificant predictor of transitivity violations ($z=-1.13$, $p=.25$). Finally, removing GSF from the full model did not reduce model fit ($BIC_{Full\ Model}=2314$, $BIC_{Full\ Model-GSF}=2222$). Similar to the change of mind analysis, difference in dwell time did not reliably predict transitivity violations ($z=-1.11$, $p=.26$). Note that this was still true when reported confidence was excluded from the regression analysis ($z=-1.27$, $p=0.20$).

I then examined how intersubject variability in metacognitive ability affected transitivity violations. I reasoned that if a well-calibrated, explicit representation of uncertainty plays a role in guiding future decisions, participants with greater metacognitive ability will show a decrease in the number of transitivity violations when the same option was presented a second time. In line with this hypothesis participants with greater metacognitive ability showed a marked reduction in transitivity violations between the first and second presentation of the same choice ($\beta=0.85$, $SE=0.42$, $z(26)=2.03$, $p<.05$). I also confirmed that this effect was not due to a relationship between metacognition and choice instability: the total number of transitivity violations was unrelated to metacognitive accuracy ($\beta=-1.83$, $SE=1.61$, $z(26)=-1.14$, $p=0.25$).

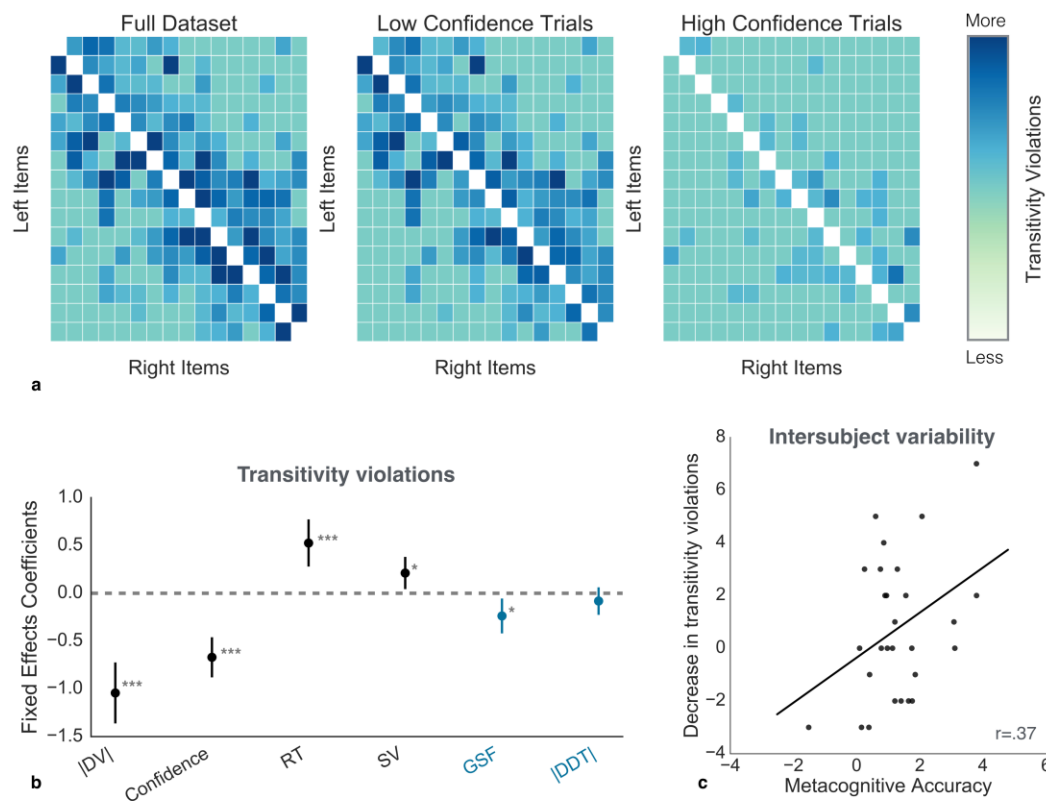


Figure 4.6. Link between Confidence and Transitivity

(a) Heat maps showing the number transitivity violations for the full sample and for high and low confidence trials (median split). The graphs are structured so that the items are ordered by increasing value from top to bottom for the items displayed on the left side of the screen and from left to right for the items displayed on the right side of the screen. The middle diagonal line is empty because no item was ever paired with itself. Note most transitivity violations took place in the low-confidence trials. **(b)** Coefficient plot for the fixed effects coefficients from a hierarchical logistic regression model predicting transitivity violations ($|DV|$ = absolute difference in value; RT = reaction time; SV = summed value; GSF = gaze shift frequency; $|DDT|$ = absolute difference in dwell time). Error bars show 95% CIs. *** = $p < .001$; ** = $p < .01$; * = $p < .05$. **(c)** Decreases in transitivity violations between the first and second presentation for each participant, as a function of metacognitive accuracy. The graph shows that participants who are more metacognitively accurate tend to become more transitive over time.

4.5. Discussion

This chapter attempted to answer two questions: What factors contribute to value-based confidence judgments and what is the benefit of explicit confidence judgments in value-based decisions? In response to the first question I found two novel factors contributing to value-based confidence judgments that predicted confidence above and beyond the well-known factors of RT and stimulus strength (Kiani et al., 2014; Kiani & Shadlen, 2009). Specifically I found that the total value of the response options predicted confidence, perhaps suggesting that some aspect of the value judgments leak into the confidence judgments, in support of the idea that the computations of confidence and value are intrinsically linked (Barron et al., 2015; Lebreton et al., 2015).

I also found that the extent to which participants switched between fixating on the options negatively predicted confidence. This might be because the brain reads these low level exploratory behaviours as one component when it constructs explicit confidence judgments. Alternatively, a low level representation of uncertainty might cause higher switch rates as a form of information seeking and also influence the post-decision confidence judgments. A third explanation would be that the BDM-values are only a noisy estimate of the true preference of the items, so the recorded unsigned difference in value ($|DV|$) is only a noisy estimate of the true unsigned difference in value ($|DV^*|$). If that is the case, both confidence and GSF might be independently driven by $|DV^*|$ and the fact that GSF predict confidence beyond recorded $|DV|$ can be explained by measurement noise. It is impossible to distinguish between these hypotheses by the results presented here, but future research could distinguish between these hypotheses behaviourally by manipulating switch-rates and testing how that influences certainty, or neurologically by finding neural substrates for GSF from EEG or MEG data (it has to be EEG and MEG rather than f(MRI) because the temporal resolution of f(MRI) is too poor to reliably link specific neural activity with GSF). These signals could then be compared to known value-and uncertainty signals.

It is worth noting that GSF is different from the more established eye metric DDT. Conceptually and empirically the two are only weakly correlated, and behaviourally they seem to have distinct effects. My work supports Cavanagh's (2014) finding that DDT has an independent, additive effect to difference in value. Or, in other words, the longer you look at an item the more likely you are to choose that item. Note that the size of this bias towards the fixated item is independent of its value. GSF, on the other hand, appears to change the steepness of the psychometric function of the value difference, so the less the participant moves his/her

gaze back and forth between the items, the more likely he/she is to choose the highest value item, in line with the idea that GSF is a low level measure of uncertainty. However, GSF did not seem to correlate with choice in the two armed bandit study.

There are a number of potential explanations for this discrepancy. First, maybe the stimuli need to have some hedonic value to trigger the kind of eye behaviour we see in the two snack experiments. Second, maybe GSF only relates to uncertainty when the stimuli are visually complex. Third, maybe GSF only works when the participants are forced to sample one item at a time. The first explanation is unlikely because we observe the influence of dwell time on choice in the bandit task and it seems unlikely that the diagnosticity of GSF is tied to the hedonic value of the items, but the dwell time effect on choice is not. Additionally we replicated the effect of dwell time and GSF for perceptual choices (see the next chapter), and those perceptual stimuli lacked hedonic value to the same extent as the bandits. The second explanation is more plausible as the visual stimuli in the snack task and perceptual task were more complex than those in the bandit task. The third explanation, that GSF only captures uncertainty when participants have to sample one item at a time is also congruent with the data. The perceptual task and the second snack experiment only showed the item participants fixated on, meaning that they could only take in visual information from one item at the time; the bandits, on the other hand, were comparably close together on the screen and their presentation was not gaze contingent, meaning that participants might have attended to a different item than the one they looked at, thus reducing the diagnosticity of GSF. This explanation is somewhat contradicted by the fact that dwell time still influenced choice in the two-armed bandit task and that the first experiment was not gaze contingent and still showed the GSF effect. Clearly further work is needed to establish the boundary conditions for when GSF captures decision-uncertainty and when it does not. If these conditions could be better understood GSF could be utilised as an easy way to capture trial-by-trial uncertainty in non-human animals (captured by head movements in rodents). This would be a valuable methodological contribution as previous work in animals has measured uncertainty as post-decision wagering (Kepecs & Mainen, 2012; Lak et al., 2014). The downside of post-decision wagering is that it behaviourally combines the choices the animals make with their confidence judgments, and conflates uncertainty estimates with economic preferences such as loss aversion (Fleming & Dolan, 2010).

This brings us to the second question: What is the benefit of having an explicit confidence representation in value-based choices? Previous work has shown that tracking decision uncertainty can improve learning (Meyniel, Schlunegger, et al., 2015), help agents to determine whether to exploit the current option or explore alternatives (Badre, Doll, Long, & Frank, 2012;

Daw, O’doherly, Dayan, Seymour, & Dolan, 2006), or allow individuals to evaluate evidence in favour of alternatives (Boorman et al., 2009). Here I found that low confidence in a choice increased the chances of changes of mind when the same option(s) appeared again. Additionally I found that the relationship between low confidence in the first presentation and subsequent changes of mind was stronger for individuals with higher meta-cognitive accuracy (whose confidence judgments better reflected whether they had chosen the items they valued more). At the neural level, such a link between meta-cognitive accuracy on the one hand and confidence-driven behaviours on the other are not entirely surprising. The rostrolateral prefrontal cortex has been shown to be central both in tracking trial-by-trial variations in confidence (Rushworth, Kolling, Sallet, & Mars, 2012; Yoshida & Ishii, 2006) and modulating uncertainty-driven behaviours (Badre et al., 2012; Boorman et al., 2009; Daw et al., 2006; Meyniel, Schlunegger, & Dehaene, 2015; Payzan-LeNestour, Dunne, Bossaerts, & O’Doherty, 2013). Additionally the rostrolateral prefrontal cortex and the frontal poles have been associated with metacognitive abilities in a number of studies (De Martino et al., 2013; Fleming et al., 2010; Lau & Rosenthal, 2011; Rounis et al., 2010). In other words, the relationship between meta-cognitive accuracy and changes of mind reported here fits a neurocomputational model where the same brain network is responsible for estimating trial-by-trial uncertainty and using that uncertainty estimate to guide behaviour, even though this study did not test neural hypotheses directly.

It is important to point out that the observed relationship between confidence and changes of mind can also be explained without postulating a causal relationship between the two. As mentioned earlier, recorded $|DV|$ might just be a noisy measure of $|DV^*|$. As such, confidence low confidence might just signal that $|DV^*|$ is smaller than is implied by our measure $|DV|$, therefore confidence would simply diagnose trials where the decision process is particularly noisy, and is subsequently more likely to be associated with a different choice when the options repeat. According to this account, confidence provides additional information to the experimenter about when a participant is likely to change their mind, but both the change of mind and the confidence judgement is a consequence of the internal computation leading up to choice, confidence has no special causal influence over changes of mind. This account cannot be ruled out from the current data, but an experimental design that influences confidence while keeping first-order performance constant might provide some insight (there are a number of such designs in relation to metamemory, see Chapter 1). If trials with artificially modulated confidence show a greater or lower incidence of changes of mind relative to control trials despite choice accuracy being the same would provide support for the causal role of confidence, if on

the other hand, the confidence manipulation did not influence subsequent changes of mind that would damage the credibility of the causal role of confidence.

So do changes of mind lead to better decisions? It is problematic to talk about correct and incorrect choices in the value domain because the value of a choice is inherently subjective, but internal consistency (transitivity) might serve as a proxy for accuracy in the value domain. I used a novel mathematical algorithm (Pedings et al., 2012) to tag choices that violated transitivity and showed that high-confidence decisions are more likely to be transitive, and hence more rational. This suggests a potential function of confidence in value-based decision making, because low confidence highlights choices that are more likely to violate the agent's overall preference patterns, and thus leave them open for economic exploitation. Furthermore, individuals who have higher metacognitive performance showed a greater reduction in transitivity violations, suggesting that people who have a more accurate internal representation of uncertainty are also able to use that representation to improve their decisions. It seems likely that the role of confidence in value-based judgments is closely related to the relationship between confidence and error correction in perceptual judgments as suggested by Boldt & Yeung (2015) and Yeung & Summerfield (2014, 2012) and recently confirmed in our lab (Kaanders et al., in preparation). This relationship with perceptual decision making highlights another strength of this study: applying psychophysics paradigms to value-based decision making (Summerfield & Tsetsos, 2012). By finding ways to parametrically manipulate value and repeat a great number of trials inside the same people, researchers can ask more precise questions about the processes informing value-based decisions than has historically been the norm (Houser & McCabe, 2014).

To summarise, these studies have extended our understanding of confidence in value-based decision making by finding that the summed value of the options and the amount of gaze shifts between options predicts confidence above and beyond value difference and response time. I have showed that confidence judgments after an initial decision predict the probability that the agents will change their minds in subsequent decisions with the same options. I have also shown that the relationship between low initial confidence judgments and changes of mind is stronger for highly metacognitive individuals. Lastly I have shown that confidence tends to be higher for transitive choices. This may suggest that low confidence highlights choices that might be irrational to the decision-maker, allowing them to change their mind and improve their choices over time.

5. The Timing of Confidence Judgments Influences Metacognitive Performance

5.1. Summary

It has been reported that the accuracy of confidence judgments differs as a function of when they are elicited in relation to the choice they evaluate. However, to the best of my knowledge, no study has compared these two confidence timings within the same individuals. Applying a within-participant comparison where only the timing of the confidence judgment relative to choice differs between conditions, I find that simultaneous confidence judgments are associated with lower second-order performance despite similar first-order performance. I investigated whether this difference in metacognitive efficiency can be attributed to differences in processing time with inconclusive results. Sequential confidence judgments were more sensitive to response time than simultaneous confidence judgments. A mediation analysis suggests that this greater sensitivity to response time fully explains the higher metacognitive accuracy of the sequential judgments. Additionally, the influence of dwell time on choice previously reported in the value domain replicates in the perceptual domain, and the association between gaze-shift frequency and confidence that was reported in the previous chapter replicates as well.

5.2. Introduction

An enigma lies at the heart of the study of confidence. On the one hand, almost every human has an intuitive sense of confidence, and if we are asked to give confidence judgments for our decisions, these judgments tend to track decision accuracy (Fleming & Lau, 2014; Nelson & Narens, 1994). On the other hand, there is an on-going controversy among the world's leading experts about what this subjective sense of confidence actually corresponds to (Fleming and Daw, 2017). Specifically, there are disagreements about two central aspects of confidence: First, what is the underlying computational quantity that confidence captures, i.e. is confidence a measure of the probability that a choice is correct or some heuristic approximation? (Aitchison, Bang, Bahrami, & Latham, 2015; Navajas et al., 2017) Second, how is confidence computed; is it based on the same internal evidence computation as the choice it evaluates? (De Martino et al., 2013; Fleming & Daw, 2017; Kiani et al., 2014; Maniscalco & Lau, 2016; Van Den Berg et al., 2016). It is possible that the conflicting findings fuelling these controversies may partially be explained if the nature of the confidence computation depends on the task structure in which it

is evaluated. Yeung and Summerfield (2012) expressed a similar idea. They pointed out that error correction and error monitoring, though long thought to have the same substrate, seem to rely on different computations. Because confidence judgments seem to be driven by the same process as error monitoring (Boldt & Yeung, 2015), it is possible that different forms of confidence judgments depend on different computations. Fleming and Daw (2017) approached this idea from a different angle: in their recent review they suggest a signal detection theoretic model that can account for all the reported findings by modelling choices and confidence judgments as resulting from separate but correlated evidence streams. By varying the strength of the evidence correlation and the relative level of noise in each evidence stream, they can account for all of the major confidence findings with the assumption that these variables vary as the result of task structure.

One simple variation in task structure that seems to influence the computation of confidence is the timing of the confidence judgment relative to the choice. To the best of my knowledge, the first study to empirically test how confidence differed for different response timings was conducted by Aitchison, Bang, Bahrami, and Latham (2015). Aitchinson and colleagues (2015) compared two groups of people who both performed the same visual discrimination task, where one group reported confidence sequentially and the other group reported their choices and their confidence judgments simultaneously. The confidence judgments of the first group were a monotonic function of a Bayes-optimal computation of the probability of being correct ($p(\text{correct})$), given the current state of evidence, whereas the confidence judgments of the second group was a mixture of $p(\text{correct})$ and a heuristic approximation (the stimulus strength of the chosen option). More recently, Siedlecka and colleagues (2016) found that metacognitive accuracy (the extent to which confidence was diagnostic of performance) was higher when participants reported their decision before their confidence, rather than the reverse. Participants were presented with masked scrambled letters and had to judge if a subsequently shown word was an anagram of the scrambled letters, with the order of the choice and the confidence rating reversed between conditions (it also contained a third condition when confidence was given before the target word was shown, but this is irrelevant for this discussion and will be ignored). Fleming and Daw (2017) explain this performance difference by arguing that the choice itself can influence the confidence judgment when the internal evidence streams driving confidence and choice are weakly correlated (when they are strongly correlated the same information influences the choice and the confidence, so the confidence judgment does not receive any new information from observing the choice). In other words, Fleming and Daw argue that confidence judgments following choice might be computationally distinct from confidence

judgments preceding them. Researchers who believe that a single evidence stream drives both confidence judgments and choices might explain the superior metacognitive performance of the retrospective confidence judgments by pointing out that participants had more time to process the information relative to the prospective confidence judgments (Van Den Berg et al., 2016). However, that explanation is unlikely in this case because first order accuracy was similar across conditions, and there is no reason why additional processing time should influence first-order and second-order performance differently if they both rely on the same evidence stream.

Finally, Kiani, Corthell and Shadlen (2014) reported discrepant results in the patterns between confidence, stimulus strength and accuracy for confidence judgments reported simultaneously with choices relative to the pattern that is commonly found for sequential confidence judgments. Typically confidence is positively associated with stimulus strength for correct decisions but show a negative or flat relationship with confidence during error trials (Kepecs et al., 2008; Lak et al., 2014; Sanders et al., 2016; Chapter 3 in this dissertation). Kiani and colleagues found a positive relationship between confidence and stimulus strength for both correct and error trials when the confidence judgment was reported simultaneously with the choice. They theorised that the reason for this pattern was that simultaneous confidence judgments force the participants to use the same evidence for correct and error trials, so stimulus strength has the same influence on both trial types, whereas sequential confidence judgments allow for additional post-choice processing which enables error detection. This error detection is stronger for trials with higher stimulus strength because the additional evidence samples from the extra processing time are more likely to contradict the error. Fleming and Daw (2017) used simulations to show that this account is untrue: even in cases when confidence and choice derive from exactly the same evidence signal, stimulus strength would be positively associated with confidence for correct trials but negatively associated with confidence for error trials. The reason for this is technical but relates to the dissociation between the internal evidence signal (that is known to the participant but unknown to the experimenter) and the stimulus strength (which is known by the experimenter, but not by the participant). However, recently Navajas and colleagues (2017) provided a new explanation for Kiani's findings. They demonstrated that confidence judgments show the traditional pattern when they computationally capture the $p(\text{correct})$, but that they show the pattern reported by Kiani et al. when they computationally capture the precision of the internal evidence. Together, these findings imply that confidence might be computed differently when it is reported simultaneously with choice, relative to when it is reported following choice.

The studies reviewed above provide an indication that confidence may differ systematically depending on when it is timed relative to choice, but the existing evidence is not conclusive. The

main reason for this uncertainty is that the studies that have tested the effect of timing on confidence judgments have used between-participants designs. Aitchinson and colleagues (2015), who found systematic differences with regard to what computational quantity best captured the confidence judgments, compared 15 participants who provided simultaneous judgments with 11 participants who provided sequential judgments. This is a major problem, because Navajas et al. (2017) have since established that there are individual differences in the computational underpinnings of confidence that are independent of task structure. Siedlecka and colleagues did not report how many participants they had in each group, but they analysed data from a total of 86 participants split between three conditions. Provided that these participants were divided evenly, that translates to approximately 29 participants per group; given the variation in metacognitive accuracy between individuals (Ais et al., 2016; Song et al., 2011), this is a fairly small sample. Lastly, while Kiani et al., (2014) did have four participants complete both simultaneous and sequential confidence reporting in the same task, the effect of these conditions was not formally compared in the paper, and only some of the data from the sequential trials were reported in the supplementary materials. Therefore, to better explore the systematic differences between simultaneous and sequential confidence judgments, a within-participant design is desirable.

Assuming that metacognitive accuracy does differ systematically between sequentially and simultaneously reported confidence judgments, it would be interesting to know what might cause such a difference. The previous chapter showed that confidence was a function of stimulus strength, response time, and eye behaviours in the value domain. In this study I recorded the same variables to test if they influence sequentially and simultaneously reported confidence differently. Additionally, recording eye movements here allowed me to test if eye behaviours inform choice and confidence in a similar way in the perceptual domain as they do in the value domain. Specifically, I wanted to test whether the time spent looking at an option influences whether it is chosen, as it does in the value domain (Cavanagh et al., 2014; Folke et al., 2016; Krajbich et al., 2010; Krajbich & Rangel, 2011; Lim et al., 2011). I also wanted to test whether the number of times participants shifted between looking at each option negatively predicted confidence and choice accuracy above and beyond stimulus strength and response time, in line with my results in the value domain (Folke et al., 2016).

To summarise, previous research suggests that the timing of confidence judgments relative to choice may influence various aspects of the confidence computation. I want to extend this work by comparing confidence judgments that are reported simultaneously with the choice to

confidence judgments that are reported sequentially after a choice. Specifically, I want to answer the following questions:

1. Siedlecka and colleagues (2016) found that participants who made their confidence judgments after they reported their choice were more metacognitively accurate than participants who made their confidence judgments before they reported their choice. Can I replicate this timing effect when the same participants complete both conditions?
2. If I discover a difference in metacognitive accuracy, can it be explained by differences in processing time between the conditions, or does it imply that the timing conditions somehow influence the computations that underlie confidence?
3. I established a set of predictors of confidence in the value domain (Folke et al., 2016); do these predictors influence confidence differently depending on whether it is reported sequentially or simultaneously?
4. If they do, can these differences explain any observed difference in metacognitive accuracy?
5. Does confidence relate to stimulus strength differently for correct and error trials when confidence judgments are reported simultaneously relative to when they are reported sequentially? If the pattern differs, it might be indicative of distinct computational underpinnings for the two types of confidence judgments (Navajas et al., 2017).
6. Does eye behaviour inform choice and confidence in perceptual decision making as it does in value-based decision making?

In order to answer these questions, I revisited the task from Chapter 3, where participants had to pick which one of two circles contained more dots, and the dot difference changed according to a 1-up-2-down staircase procedure. This procedure has many benefits. A perceptual two-alternative forced choice task is the bread and butter of confidence research, so using such a design assures comparability with much of the existing literature. The fact that participants have to compare the stimulus strength of two stimuli means that they have to shift their visual attention between two different objects. Thus, I can test whether the relationships between eye behaviour and choice that I explored in the value domain work in the same way for perceptual discrimination. Finally, the staircase procedure keeps first order accuracy fixed between sessions and participants, so any observed variations in metacognition will not be contaminated by variation in first-order performance. The task was also slightly modified from Chapter 3, in that I made the presentation of the dots gaze-contingent to ensure that participants could only sample information from one source at the time. Additionally, participants used fixations to report their choices and their confidence judgments. This was done so that confidence and choice would

result from a single ballistic motor action for the simultaneous condition to ensure that choice and confidence was determined simultaneously and to minimise the motor latency between when the decision was made and when the response was recorded (Kiani et al., 2014). Confidence and choices were also reported by fixations in the sequential condition to make the conditions as similar as possible. For these reasons, this experimental design is well suited for answering the research questions and exploring the influence of the timing of confidence judgments on metacognition.

5.3. Methods

5.3.1. Experimental Procedure

Participants completed a two-alternative forced choice visual discrimination task. They were instructed to determine which one of two circles on a computer screen contained more dots. The presentation of the dots was gaze-contingent, so that participants only saw the dots of a given circle when they looked inside it. The participants could look inside either circle however many times they wanted for however long they wanted. Participants were instructed to respond as fast as they could while maintaining accuracy, but apart from this verbal instruction, response times were not constrained. The difference in dots was updated online in a 1-up-2-down staircase procedure to keep performance close to 70% correct. Participants indicated their confidence in their choice on a visual scale. The scale was unlabelled but colour coded with a transition between blue and green (see Figure 5.1.). Participants were instructed to treat the blue edge as “guessing” and the green edge as “certain”. First-order and second-order performance in this task was incentivised according to a linear scoring rule described later in the methods. The experimental trials were split into two sessions that were conducted on different days. In one session participants gave their confidence at the same time they made their choice. In the other session, they gave their confidence after they had reported their choice. The order of the sessions was counterbalanced between participants. Participants indicated their choices and confidence judgments with their eyes (see Figure 5.1.). This ensured that participants could provide their confidence judgments and choices simultaneously with a single ballistic eye movement, in line with previous work (see Kiani & Shadlen, 2014).

Each session consisted of 300 trials split into 12 blocks with 25 trials in each. The first two blocks in each session consisted of practice trials that were used to calibrate the staircase procedure that determined the difficulty of the task. Participants were informed that their performance on these trials would not count towards their earnings and that these trials were

excluded from all analyses. Participants received feedback after their choices during the first block in each session. During the practice blocks the experimenter monitored how the participants used the confidence scale, and reminded them of the instructions if they consistently provided confidence judgements near the extremes of the scale. For the remaining blocks participants received no feedback on their performance apart from seeing their cumulative total score between blocks.

Performance during the main blocks was incentivised according to a linear scoring rule. Points = $50 + \text{correct} \times (\text{confidence} \div 2)$ (see Figure 5.1. c), where correct was a dummy coded as 1 for correct trials and as -1 for error trials, and confidence was scaled from 0 (guessing) to 100 (certain). In other words, participants earned more points for being highly confident when they were correct and less points for being highly confident when they made errors. Participants were informed that every 4000 points corresponded to an extra pound in their participation payment (rounded to the nearest ten cents at the end of the session) and were shown their total score between blocks in order to maintain their motivation throughout the task. The reason why there was no option to report “certain errors” in the sequential condition is two-fold: first, I wanted to maintain the same scale for both conditions to ensure comparability, second when piloting an untimed version of this perceptual discrimination task I included an option to report certain (motor) errors and it was never used, so I judged it unlikely that excluding it would meaningfully impact the results. The gaze-dependent dot-discrimination tasks were programmed in Experiment Builder version 1.10.1640, SR-Research, Ontario.

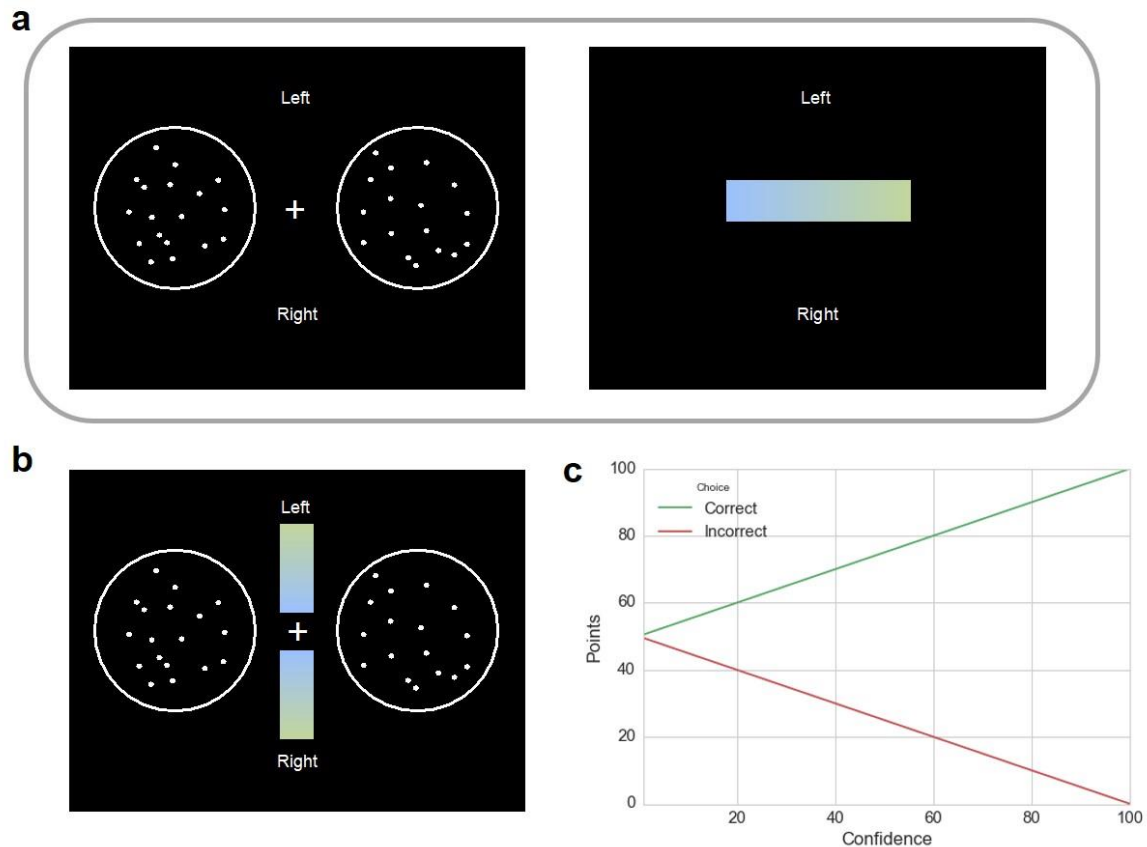


Figure 5.1. Experimental Procedure

Participants completed two sessions of a forced choice perceptual discrimination task, which involved picking the circle containing the most dots. The presentation of the dots was gaze-contingent so participants only saw the dots inside the circle they were fixating on. During one session **(a)** they first made their choice by fixating on the words “left” or “right” on the top and bottom of the screen and then gave a confidence judgment by fixating somewhere on a visual scale in the middle of the screen. During the other session **(b)** they made their choice and confidence judgment simultaneously by fixating somewhere on the upper confidence scale for the left choice and somewhere on confidence scale below for the right choice. The size of the confidence scales was the same across sessions. The order of the two sessions, whether left was associated with up or down and whether high confidence were to the left or the right in the two step session were all counterbalanced between participants. Participants were instructed to respond as quickly as they could while still being accurate, but other than that their response times and interjudgment times were unconstrained. **(c)** The choices and confidence judgments were incentivised according to a linear scoring rule, so that they earned more points for high confidence judgments when they were correct and less points for high confidence judgments when they were incorrect.

5.3.2. Participants

32 participants (25 female) completed this study (an additional 3 started the experiment but had to abort because of technical difficulties with the eye tracker). The mean age of the sample was 26.3 years (sd = 6.6). All participants gave informed consent prior to participating. They were paid a £10 show-up fee for each session they attended and could earn up to £10 more based on their performance. Actual payments ranged between: £26.60 and £28.50.

5.3.3. Eye Tracking

Eye movements were recorded at 1000Hz with an EyeLink 1000 Plus eye tracker (SR-Research, Ontario). Participants sat approximately 60 cm away from a 68.58 by 59.77 cm computer screen. A headrest kept their heads immobile during the blocks to maximise the accuracy of the eye tracker. Participants were told to report if their confidence judgments did not land where intended, in which case the eye tracker was recalibrated. *Areas of interest (AIs)* were predefined by the 2 squares with the same length as the circumference of the stimulus circles containing the dots (400 x 400 pixels). I derived two variables from the eye tracking data: the total dwell time in each AI for a given trial, and gaze shift frequency (GSF). GSF measured the number of times participants shifted immediately between the AIs (saccades that started in one AI and ended in the other AI). Saccades within the same AI and to other parts of the screen were ignored.

5.3.4. Hierarchical Models

Note that all predictors entered into the hierarchical models are z-scored on the participant level, and that response time was log-transformed prior to being z-scored to make RT distributions approximately normal.

4.3.5. Drift Diffusion Models

I fitted three hierarchical drift diffusion models (DDM) to the combined data from both sessions to evaluate whether fixation time influenced evidence accumulation in an additive or interactive fashion (see Chapter 2 for details on the software and fitting procedure). The three models are similar to the models reported in Chapter 4 with the difference that they predicted picking the circle with most dots rather than the highest valued snack item and that they used the number of dots in each circle rather than the BDM value of the options as a measure of stimulus strength. The interactive model had not converged after 2000 samples (Gelman-Rubin>1.1), so I reran the analysis with 4000 samples per chain and a burn-in of 200. The longer chains showed good mixing, indicating convergence (Gelman-Rubin=1.02).

5.4. Results

5.4.1. First Order Choices

Before evaluating whether metacognitive accuracy changed as a function of the timing of the confidence judgments, I wanted to compare whether the first order performance differed between the sessions. As Table 5.1. indicates, accuracy, stimulus strength (dot difference), response times and total dwell time all seem to be similar across experiments. These similarities are illustrated by the Bayes Factors which suggest weak evidence in favour of the hypothesis that these quantities are being drawn from the same distribution in both sessions.

Table 5.1. First-order response comparisons (DF=31)

	Simultaneous	Sequential	t	p	BF
Accuracy (%)	72.1 (1.1)	72.2 (1.5)	0.55	0.59	4.6 in favour of H_0
Dot Difference (dots)	3.36 (0.90)	3.54 (0.86)	1.61	0.11	1.7 in favour of H_0
RT (ms)	4283 (1969)	4065 (1470)	1.01	0.32	3.3 in favour of H_0
Total Dwell Time (ms)	2987 (1669)	2964(1216)	0.13	0.90	5.3 in favour of H_0
GSF	3.59 (1.51)	3.80 (1.13)	1.29	0.21	2.5 In favour of H_0

While exploring the first order responses I found a response bias favouring the right option over the left option. For example, in a hierarchical logistic model predicting correct responses from stimulus strength, response timing condition, and an interaction term between the two, the intercept was shifted to favour the right option ($z=2.76$, $p<.01$) and this effect was stronger for the sequential trials than the simultaneous trials ($z=3.80$, $p<.001$). I have not been able to find the cause of this response bias, but it should balance out as the position of the correct options was randomised (left was the correct option 49% of the time for both response timings).

However, to ensure that this spatial bias does not unduly influence my findings, I will control for it algebraically by including it in all regression models that predict first-order accuracy or confidence.

I also wanted to test whether dwell times influenced choices additively, as they did for value (see Chapter 4). As for Chapter 4, I compared three DDM models: a null model that predicted drift rate from dot difference, an additive model that predicted drift rate from the additive effect between dot difference and dwell time ratios, and an interactive model that weighted the dots in each circle by how long that circle was fixated. In line with previous research, the additive model outperformed the other two models ($DIC_0=79566$, $DIC_a=78937$, $DIC_i=79169$). The winning

additive model showed that the influences of dot difference and difference in dwell ratio on choice were both positive ($\text{coefficient}_{\text{dot difference}}=0.07(0.003)$, $\text{coefficient}_{\text{dwell ratio}}=0.81(0.03)$).

5.4.2. Three-way Interaction Between Confidence, Stimulus Strength and Accuracy does not Depend on Timing of Confidence Judgments

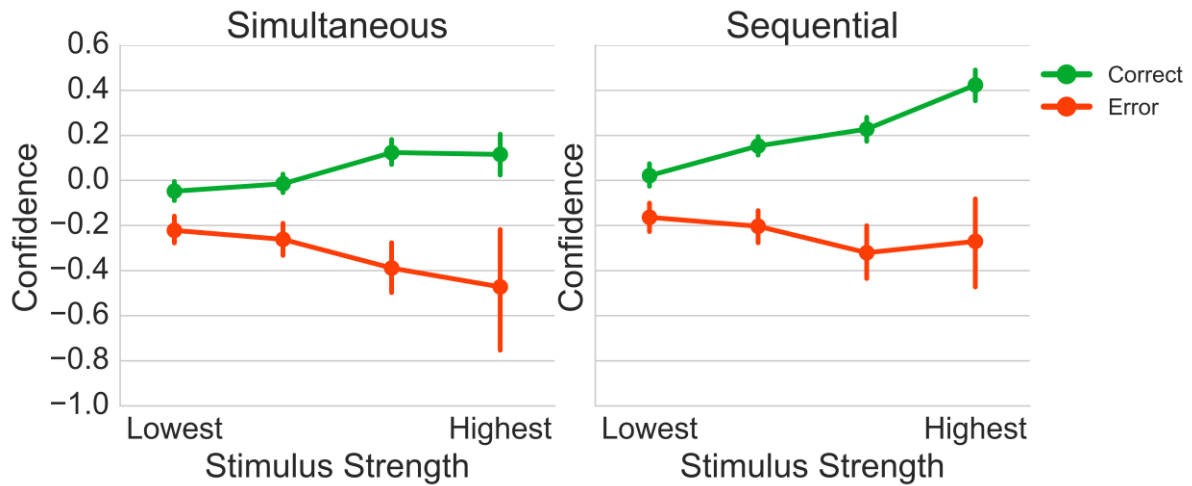


Figure 5.2. The Relationship Between Confidence and Stimulus Strength for Simultaneous and Sequential Responses.

Stimulus strength corresponds to dot difference, quartile binned on the participant level. The smallest dot difference is furthest to the left, the greatest dot difference is furthest to the right. Confidence was z-scored at the participant level. Confidence was typically higher for the sequential confidence judgements relative to the simultaneous confidence judgements. For both the sequential and simultaneous confidence judgements confidence increased with stimulus strength for the correct trials and decreased with stimulus strength for error trials.

In order to explore how confidence related to stimulus strength for correct and error trials, I ran two hierarchical linear regression models with participantwise variation in intercepts and fixed slopes, one model focusing on error trials and one model focusing on correct trials. During the correct trials sequential confidence judgments were in general higher than simultaneous confidence judgments ($t=5.35$ $p<10^{-7}$; see Figure 5.2), stimulus strength was associated with more confident responses for both conditions ($t=5.85$ $p<10^{-8}$) but this effect was stronger for the sequential confidence judgments ($t=3.05$ $p<.01$). During the error trials, there was no difference in confidence between the conditions when the role of stimulus strength was controlled for ($t=0.46$, $p=.65$), stimulus strength negatively predicted confidence ($t=-3.36$, $p<.001$) to the same extent for simultaneous and sequential trials ($t=0.59$, $p=.56$). Additionally, the spatial bias in choice persisted in confidence in that participants were less confident when the

left option was chosen when the dot difference of the trial was controlled for both for error trials ($t=-4.26$, $p>10^{-4}$) and correct trials ($t=-8.49$, $p>10^{-10}$).

5.4.3. Metacognitive Efficiency is Higher for Sequentially Reported Confidence than Simultaneously Reported Confidence

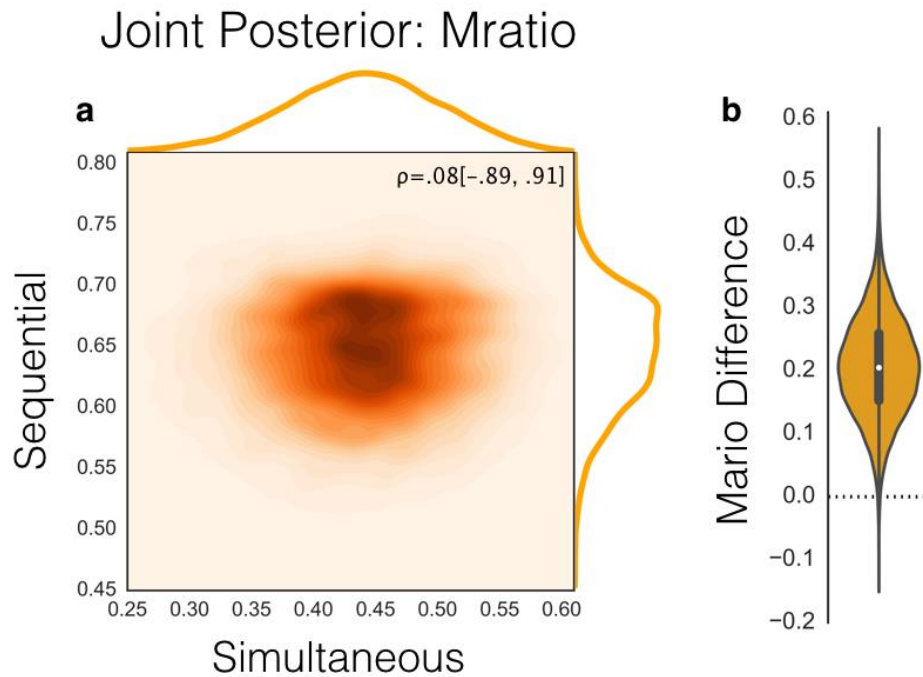


Figure 5.3. Sequential Confidence Responses are Associated with Higher Metacognitive Efficiency

(a) The joint posterior distribution of the mean Mratio for the sequential trials (Y) and the Simultaneous trials (X). ρ in the upper right corner reports the mean and the 95% HDI for spearman's rho between the parameters. There does not appear to be a strong correlation between metacognitive performance in the sequential and simultaneous versions of the dot discrimination task. **(b)** The mean Mratio difference between sequential confidence ratings and simultaneous ratings as estimated by the hierarchical Bayesian model. The Mratios are higher for the sequential confidence judgements than the simultaneous confidence judgments.

To estimate metacognitive efficiency I used the Mratio, a ratio of metacognitive sensitivity, captured by meta- d' , and first-order sensitivity, captured by d' . I used a hierarchical Bayesian estimation method discussed in Chapter 2 to fit the Mratios (Fleming, 2017). The confidence data were binned into participantwise tertiles prior to the Mratio estimation. I fitted both

conditions simultaneously by estimating the mean Mratio for each condition from a bivariate distribution (Figure 5.3. a). I did this to account for the paired nature of the estimate, as each participant would have an estimate of metacognitive efficiency from each of the timing conditions. Two results are noteworthy: first, the posteriors for the sequential and the simultaneous confidence responses are hardly correlated ($\rho=.08$, $\text{HDI}=(-.89, .91)$), but note that this might just reflect the extreme uncertainty around this estimate. Second, the posterior mean Mratio for the sequential confidence judgments was much higher than for the simultaneous confidence judgments. I subtracted the posterior mean of the sequential confidence judgments with the posterior mean from the simultaneous confidence judgments to quantify this difference. 99.56% of the probability mass was greater than 0 (posterior odds =228 :1; Figure 5.3. b), providing strong evidence that sequentially reported confidence resulted in higher metacognitive efficiency than simultaneously reported confidence. Note that this difference cannot be explained by differences in the distribution of confidence responses between the sessions. Mean confidence judgments were slightly (but not significantly) higher for the sequential sessions compared to the simultaneous sessions (mean difference = 3.97 points on a 100-point scale, $t(31)=1.84$, $p=.08$, $\text{BF}=1.2$ in favour of H_0). The participantwise standard deviations were similar across sessions (mean difference = 0.12 points on a 100-point scale, $t(31)=0.13$, $p=.90$, $\text{BF}=5.3$ in favour of H_0). Mean confidence judgments were highly correlated across the conditions ($r=.80$, $t(30)=7.39$, $p<10^{-7}$), as were the standard deviations ($r=.64$, $t(30)=4.52$, $p<10^{-4}$). Appendix 2 shows the confidence distributions for each participant.

5.4.4. Can Differences in Processing Time Account for Differences in Metacognitive

Efficiency?

The previous section established that the marked difference in metacognitive efficiency between sequential and simultaneous confidence judgments could not be explained by any obvious differences in the distribution of confidence judgments. Next, I wanted to test whether the difference in metacognitive efficiency could be explained by differences in processing time. While the simultaneous and sequential trials had similar response times (see Table 5.1.), the sequential trials had additional interjudgment time (IJT) between the choice and the confidence judgment, so the total response times ($\text{RT} + \text{IJT}$) were longer for the sequential sessions (mean difference = 705 ms, $t(31)=3.19$, $p=.003$, $\text{BF}=11.6$ in favour of H_1). While variations in IJT were unrelated to variations in confidence, difficulty, or first-order or second-order accuracy (see Appendix 3) it is still possible that the extra processing time they afforded caused the observed increase in second-order performance.

If the within-participant difference in metacognitive accuracy between the sequential and simultaneous trials could be explained by additional processing time, the participants who showed an improvement in metacognitive efficiency between sessions should also have slower total response times in the sequential session relative to the simultaneous session. To test this hypothesis I ran a correlation comparing the difference in M_{ratios} between the sequential session and the simultaneous session with the difference of mean total response times between the sessions (see Figure 5.4. c). I found that there was no relationship between metacognitive gains for the sequential session and increases in total response time ($r = -.01$, $t(30) = -0.08$ $p = .94$, $BF = 3$ favouring the null).

Though variations in interjudgment times were independent of metacognitive efficiency for the sequential sessions and the increases in total RT did not predict increases in metacognitive gains, I have not ruled out that the observed difference in metacognitive efficiency is due to additional processing time. Shadlen and colleagues have shown the last ~400 ms of evidence is not utilised in the initial choice because of the lag between the internal commitment to a decision and motor onset, but that this information may drive subsequent changes of mind and confidence if participants have enough time to alter their motor response (Resulaj et al., 2009; Van Den Berg et al., 2016). It is therefore possible that metacognitive efficiency is higher for the sequential confidence judgments because they benefit from these 400 ms of extra evidence integration, while all the variation in IJT is due to variation in the motor response times rather than processing times. According to this hypothesis the sequential confidence judgments benefit from additional processing time that is constant across trials and is therefore not captured by IJT. In order to control for this possibility I wrote a sorting algorithm that matched simultaneous trials and sequential trials on total response time. For each participant, the algorithm cycled through the simultaneous trials in random order, trying to find a sequential trial with matching total RT. For each simultaneous trial, the algorithm listed the sequential trials that had a total RT ($RT + IJT$) $\pm 20\%$ of the standard deviation of simultaneous response times for that participant. If there were several possible matches the algorithm selected the sequential trial with the smallest absolute total RT difference to the simultaneous trial. If a match was found, the matching trials were stored in a separate data frame, and the sequential trial was removed from the prospect pool to prevent the same sequential trial from being matched with multiple simultaneous trials. Trials without any matches were discarded. The matched subset consisted of 7434 trials, 48% of the total sample. Each participant had an average of 116 matched trials from each session (range: 29-139). Following the matching algorithm there were no significant differences in mean trial length between the sequential sessions and the simultaneous sessions (mean difference = 0.5 ms,

$t(31)=-0.72$, $p=.48$, $BF=4.2$ in favour of H_0 ; see Figure 5.4. a). The conditions did not differ in terms of average dot difference for the matched trials (mean difference= 0.2 dots, $t(31)=-1.45$, $p=.15$, $BF=2$ in favour of H_0), but the sequential sessions were associated with slightly higher accuracy (mean difference=3%, $t(31)=2.15$, $p=.04$, $BF=1.5$ in favour of H_1).

Second-order accuracy was somewhat higher for the sequential session than the simultaneous session for the matched trials as 78.85% of the probability mass was greater than 0 (posterior odds =4:1; Figure 5.4. b). This difference is too weak to support the conclusion that Mratios are reliably higher for the sequential session when total RTs are matched.

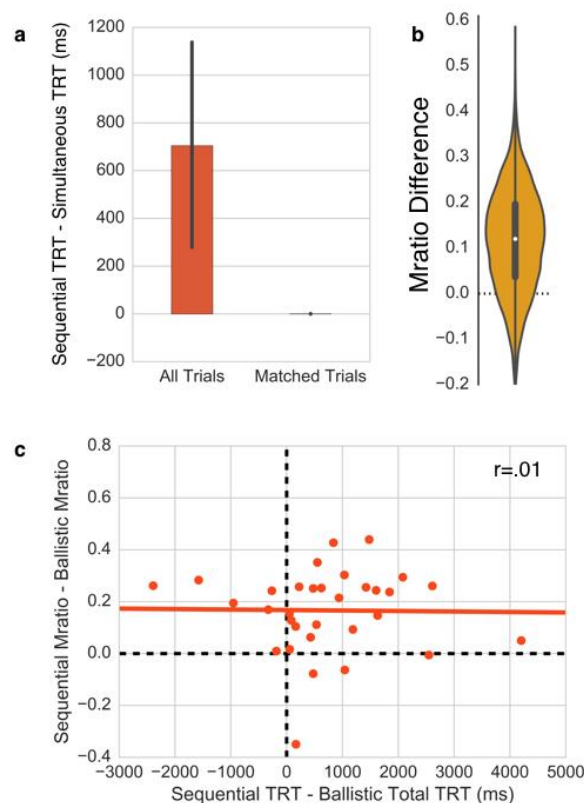


Figure 5.4. Second-order Accuracy as a Function of Total Response Time

(a) Difference in mean total response times between the sequential session and the simultaneous session for the full dataset and for a subset of trials selected to match in total response time. **(b)** The difference in Mratios between the sequential and simultaneous sessions for the subset of trials with matched response times. The sequential session is weakly but not reliably associated with higher metacognitive accuracy for the response time matched subset. **(c)** Correlation between difference in Mratios between the sessions and difference in mean total response time.

Differences in response time between the sequential and simultaneous sessions are not associated with changes in metacognitive efficiency.

5.4.5. Sequentially and Simultaneously Reported Confidence Cause Differences in Sensitivity to the Variables Predicting First-order Accuracy and Confidence

In the previous chapter I reviewed a set of variables that predicted confidence judgments and choices in the value domain. Here I tested whether the timing of the confidence responses influence how sensitive first-order accuracy and confidence judgments are to these predictors and whether any such differences can explain the observed difference in metacognitive efficiency. I will begin by exploring what variables predicted first-order accuracy.

A hierarchical logistic regression model with random slopes for the main effects but fixed effects for the interaction terms showed that stimulus strength positively predicted first-order accuracy ($z=16.33$, $p<10^{-10}$; see Figure 5.5. a) and that there was no difference in the strength of this relationship across conditions ($z=-0.91$, $p=.36$). Participants showed a spatial bias favouring the right option ($z=3.49$, $p<10^{-4}$) and this bias was stronger in the sequential trials than the simultaneous trials ($z=-3.17$, $p=.002$). Response times negatively predicted accuracy ($z=-2.66$, $p=.007$), and this effect did not differ significantly between sessions ($z=-1.50$, $p=.13$). With regards to the eye tracking variables, difference in dwell time between the correct and incorrect items predicted accuracy ($z=6.82$, $p<10^{-10}$) and this effect was stronger for the sequential sessions than the simultaneous sessions ($z=3.23$, $p=.001$; See Figure 5.5. b). Gaze shift frequency negatively predicted accuracy ($z=-3.45$, $p<10^{-4}$), in line with the idea that it is a low level signal of uncertainty (Folke et al., 2016). The strength of this effect did not differ significantly between groups but was slightly stronger for the simultaneous condition than the sequential condition ($z=1.75$, $p=.08$). Also, when all the other variables were accounted for, the intercept was still marginally higher for the sequential session ($z=2.19$, $p=.03$). Together these results show that the same variables that predicted choices in value-based decisions also influence choice in perceptual discrimination tasks. Additionally, most variables predicted accuracy to the same extent both for sequential and simultaneous trials; one noteworthy exception is difference in dwell time, which was a stronger predictor of choice for the sequential trials than the simultaneous trials.

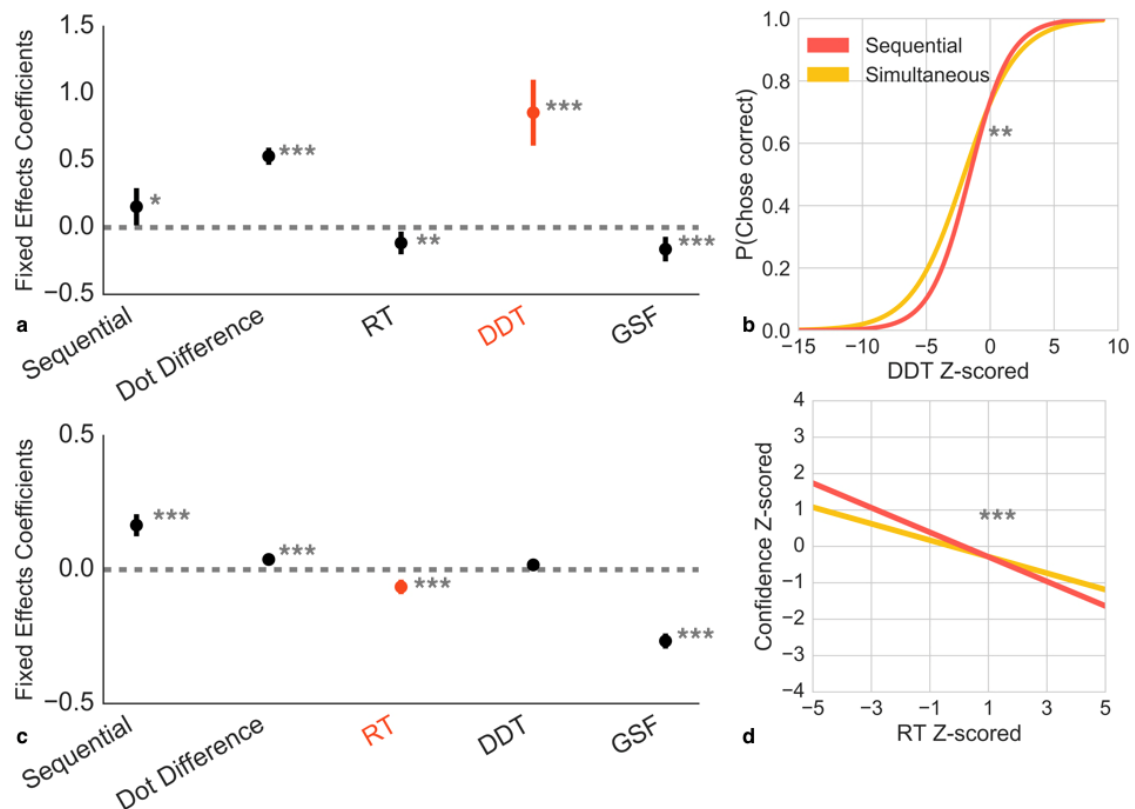


Figure 5.5. Predictors of First-order Accuracy and Confidence

(a) Coefficient plot for a hierarchical logistic regression predicting the probability of picking the correct option. DDT was the only predictor that significantly differed between the sequential and simultaneous conditions (highlighted in orange). (b) Logistic regression predicting the probability of choosing the correct option as a function of DDT (z-scored on the participant level). DDT was a stronger predictor of accuracy for the sequential trials relative to the simultaneous trials. (c) Coefficient plot for a hierarchical linear regression predicting confidence. Response time was the only predictor that significantly differed between the sequential and simultaneous conditions (highlighted in orange). (d) Linear regression predicting how response times relate to confidence (both variables z-scored on the participant level). Response time negatively predicted confidence to a greater extent in the sequential trials relative to the simultaneous trials. (RT= response time (log-transformed); DDT = difference in dwell time, GSF = gaze shift frequency). Error bars show 95% CIs. *** = $p < .001$; ** = $p < .01$; * = $p < .05$.

A hierarchical linear regression model with random intercepts but fixed slopes showed that participants were more confident for trials with larger dot differences ($t=3.56$, $p<.001$; see Figure

5.5. c) and that there was no difference in the strength of this relationship across conditions ($t=1.46$, $p=.14$). As for accuracy, there was a spatial bias favouring the right option ($t=4.53$, $p<10^{-5}$) but the strength of this confidence bias did not differ between the sessions ($t=-0.61$, $p=.54$). Faster response times predicted higher confidence ($t=-4.60$, $p<10^{-5}$), and this effect was stronger for the sequential confidence judgments ($t=-6.45$, $p<10^{-9}$; Figure 5.5. d). With regards to the eye tracking variables, difference in dwell time between the correct and incorrect items did not predict confidence judgments ($t=1.68$, $p=.09$), and there was no reliable difference in this effect between sessions ($t=0.02$, $p=.98$). Gaze shift frequency was by far the strongest predictor of confidence ($t=-18.37$, $p<10^{-10}$), and the strength of this relationship was independent of the timing of the confidence judgments ($t=1.68$, $p=.09$). Additionally, sequential confidence judgments were associated with higher confidence when all other effects were accounted for ($t=7.86$, $p<10^{-10}$). This difference was probably driven by greater confidence in correct trials (see Figure 5.2.). There are two main takeaways from this section: GSF was a strong predictor of confidence even when stimulus strength and RT was controlled for, suggesting that the relationship between GSF and confidence that I first discovered in value based choice seems to hold for perceptual decision making as well. Additionally, RT was the only variable predicting confidence that was reliably moderated by the timing of the confidence judgments. One possible reason for this difference between conditions is that RT was a cleanly and directly related to choice in the sequential version of the task whereas in the simultaneous reporting condition RT became a noisier measure of choice quality because it was contaminated by aspects of the confidence judgement (for example, previous work suggests that higher confidence judgements are associated with faster confidence RT, see Moran et al., 2015). There are two reasons to be sceptical of this account. First, if RT was a noisier measure of choice in the simultaneous trials because it was contaminated by the confidence judgement we would expect RT in the simultaneous condition to be more weakly associated with choice accuracy relative to the sequential condition, whereas RT appears to be equally diagnostic of choice accuracy in both conditions (Figure 5.5. a). Second, in the sequential condition IJT (Confidence RT) is not associated with reported confidence, choice accuracy or evidence strength. Therefore it is hard to see how confidence RT could systematically bias choice RT in the simultaneous trials as confidence RT does not seem to systematically track anything in this particular task.

Given that the sensitivity of confidence judgments to RT differed between sequential and simultaneous trials, I wanted to test whether this difference was categorical or gradual. That is, did the influence of RT increase monotonically with processing time after the choice, or was there a distinct cut-off between simultaneously and sequentially reported confidence judgments?

To compare these theories I tested whether RT interacted with IJT for the sequential trials, my reasoning being that if there was a gradual benefit of additional processing time on RT integration, higher interjudgment times would increase the sensitivity of confidence to RT. A hierarchical linear regression model with random intercepts and slopes failed to find evidence for such an interaction ($t=1.52$, $p=.14$; See Figure 5.6. a).

5.4.6. Differential Sensitivity to Response Times Explains Differences in Metacognitive Efficiency

Finally, I wanted to investigate whether the differential sensitivity of confidence to RT between the sequential and simultaneous trials could explain the difference in metacognitive sensitivity. Given that response times negatively predicted accuracy for both the sequential and simultaneous sessions and had a stronger influence on confidence during the simultaneous session, it is possible that the difference in metacognitive accuracy between the sessions can be explained by differences in how sensitive the confidence reports are to RT. To test this hypothesis I ran a hierarchical moderated mediation analysis where I predicted confidence from accuracy with RT as a mediator and dot difference and an interaction between dot difference and accuracy as covariates. I allowed intercepts and the slopes of accuracy on RT and confidence to vary between participants. Accuracy was a stronger predictor of confidence for the sequential trials than for the simultaneous trials, in line with the Mratio analyses (see Figure 5.3. b). Additionally, the RT mediation was stronger for the sequential trials than the simultaneous trials (coefficient = -0.07 , $p < 10^{-10}$) and the direct effect of accuracy on confidence did not differ between the sequential and simultaneous trials when the RT mediation was controlled for (coefficient = -0.003 , $p=.92$). In other words, it appears as if the observed difference in second-order accuracy between the simultaneous and sequential confidence judgments can be explained by the fact that confidence is more sensitive to response time in the sequential case (see Figure 5.6. b).

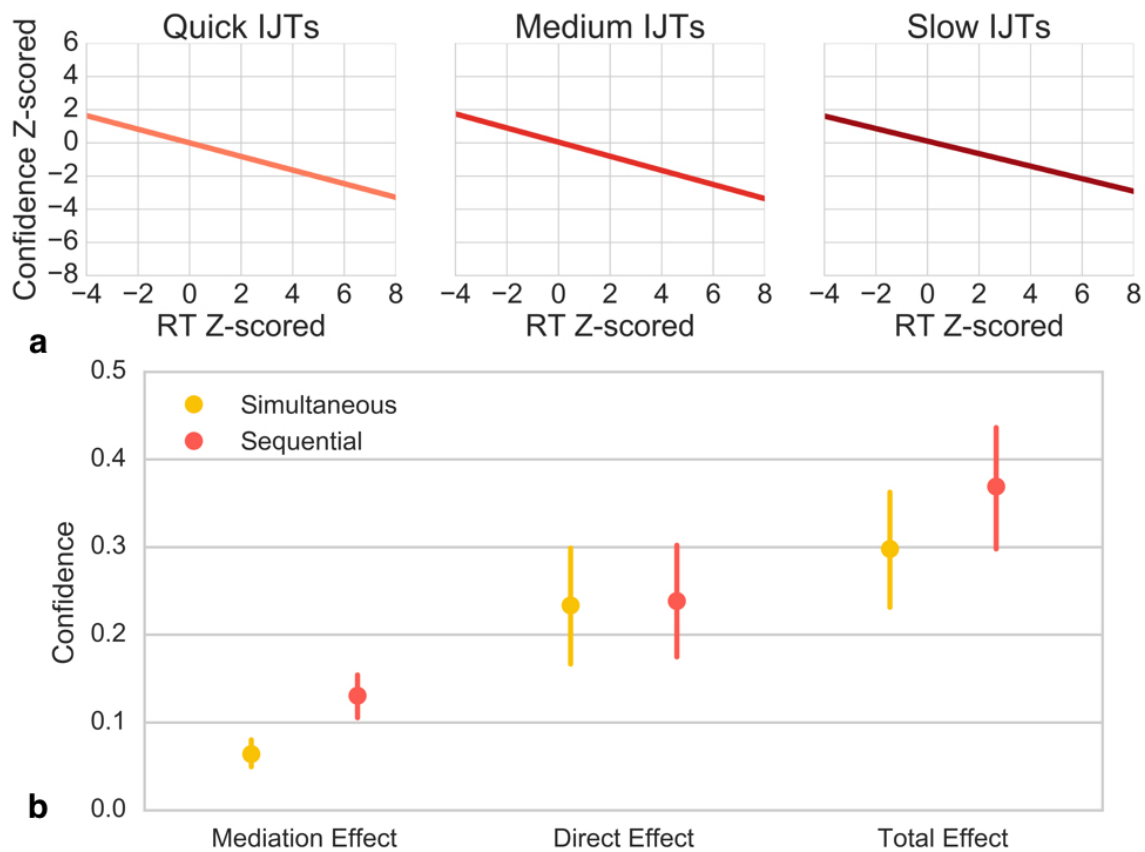


Figure 5.6. The Relationship between RT and Confidence

(a) Interjudgement time was split into participantwise tertiles; the relationship between z-scored confidence and z-scored log-transformed response times are displayed for each bin.

Interjudgment times do not mediate the influence of response time on confidence for the sequential trials. **(b)** Coefficient plot for a mediation analysis predicting confidence based on accuracy, with response condition (sequential or simultaneous) as a moderator, response time as a mediator and dot difference and an accuracy-dot difference interaction as covariates. The effect of accuracy on confidence (the direct effect) is the same across conditions when the mediation of response time is controlled for. In other words, the difference in metacognitive accuracy between the response conditions appears to be fully explained by the difference in confidence sensitivity to response time.

5.5. Discussion

This experiment compared sequentially and simultaneously reported confidence judgments in a perceptual discrimination task. First order performance was similar across the experiments with

the exception that the time spent looking at the options influenced choice more in the sequential sessions. Despite similar first-order performance in terms of difficulty, accuracy and response times across the sessions, sequential confidence judgments were associated with higher metacognitive performance than simultaneous confidence judgments. These results align well with results from previous studies on the timing of confidence judgments which used between participant designs (Aitchison et al., 2015; Siedlecka et al., 2016). Importantly, this difference could not be explained by differences in the confidence distributions between the sessions. In fact the means and standard deviations of the confidence ratings correlated highly between the sessions, in line with previous work on the stability of confidence distributions over time (Ais et al., 2016). However, in contrast to Ais and colleagues, my measure of metacognitive performance did not correlate strongly between the sessions, despite the constant task structure. This may imply that simultaneous and sequential confidence judgments depend on different abilities that are only weakly correlated, but it might also be that the current dataset was too small to reliably estimate the correlation coefficient for a multivariate normal posterior, as is suggested by the very wide confidence intervals for the correlation estimate.

There is an ongoing controversy in confidence regarding whether confidence and choices depend on a single evidence stream (see Fleming & Daw, 2017 for a review). Proponents of the single stream account would argue that the difference in metacognitive efficiency between simultaneously and sequentially reported confidence can be fully explained by sequential confidence judgments benefitting from additional processing time (Kiani et al., 2014; Van Den Berg et al., 2016). I evaluated this hypothesis in three different ways. First, I reasoned that if additional processing time leads to better confidence judgments, longer interjudgment times should be associated with better metacognitive performance for the sequential confidence judgments. In fact, longer interjudgement times have been reliably linked to better confidence discrimination within participants in a series of experiments on perceptual and knowledge-based decision making by Yu, Pleskac and Zeigenfuse (2015). However, I failed to find any evidence of this either within or between participants (see Figure 5.3, Figure 5.5 and Appendix 3). The reason for the discrepancy between my findings and those reported by Yu and colleagues is probably that the first-order judgements in their experiments were under speed stress, whereas this study had unconstrained response times. Work by Baranski and Petrusic (1998) suggests participants simultaneously process choices and confidence judgements when choice response times are unlimited, but compute confidence post-choice when response times are constrained. Second, if additional processing time drove the increase in metacognitive efficiency the participants with the greatest increase in total response time between the sequential and

simultaneous session should also see the greatest increases in metacognitive efficiency. There was no evidence of such a relationship. However, both of these tests assume that there is a monotonic relationship between metacognitive efficiency and response time. The Shadlen lab has presented evidence that the last ~400 ms of information presented before motor onset does feed into the initial choice (because the latency between an internal commitment to a decision and motor onset) but that this information may inform subsequent confidence judgements or changes of mind (Kiani & Shadlen, 2009; Van Den Berg et al., 2016). This means that confidence judgements in the sequential trials automatically benefit from ~400 ms worth of extra evidence integration. The fact that IJTs were not diagnostic of metacognitive accuracy does not count as evidence against this account as this evidence was integrated in the last 400 ms of the choice trial, so it is conceivable that variations in IJTs are mostly driven by motor-noise rather than additional processing. In order to control for this possibility I wrote a sorting algorithm that matched simultaneous trials with sequential trials of a similar length from the same participant. Note that this method slightly disfavours the sequential trials as 99% of IJTs were longer than 400 ms. First-order accuracy was higher for the sequential session than the simultaneous session for this subset of the data, but this is controlled for mathematically in the M_{ratio} estimation. Crucially the difference in difficulty between the sessions was small and not statistically significant. Though the M_{ratio} hyperparameter was still higher for sequential trials than the simultaneous trials for the matched subset of the data, this difference was not reliable enough to draw a strong conclusion. Larger sample sizes could resolve this uncertainty; a follow-up study that does not require an eye tracker but where responses could be recorded by a button press would allow for much larger samples through sites like Amazon's Mechanical Turk.

I found that DDT influenced choices in an additive fashion and that GSF influenced confidence, replicating the role of visual attention in evidence accumulation that I first discovered in the value domain in perceptual domain. A high GSF was associated both with lower-first order performance and with lower confidence when the other variables were accounted for, supporting the idea that GSF might be a low level marker of uncertainty during evidence accumulation. Difference in dwell time, on the other hand, strongly predicted first-order accuracy but did not predict confidence. This is the result of how these models were coded; because difference in dwell time was designed to predict the correctness of a response, it was constructed as the difference in dwell time between the correct option and the incorrect option, capturing the additive effect of visual attention on evidence accumulation (Cavanagh et al., 2014). However, the participant was not aware if they had focused more on the correct option or not, so they could not use that gaze behaviour as a marker of confidence. It might seem confusing that the

choice model was written to predict accurate choices in this chapter but captured the probability of picking an item in a specific reference position in the previous chapter. I chose the accuracy-based format here because it allows for a comparison of the relative predictive power of a variable for accuracy and confidence; this option was not available for the value-based experiments where no option can be said to be objectively correct.

Previous work has compared simultaneous and sequential confidence judgments in the same task (e.g. Aitchison et al., 2015; Siedlecka et al., 2016), but to the best of my knowledge this is the first within-participant comparison. Making a within-participant comparison allowed me to rule out individual differences as the cause of the observed difference in metacognitive performance between the sessions (Ais et al., 2016; Navajas et al., 2017). One example of potentially conflating effects from response timing and effects due to individual differences come from Kiani, Corthell and Shadlen (2014). They reported that simultaneous confidence responses became more positive the stronger the stimulus strength of the trial, independent of the accuracy of the first order response, in contrast to most studies that find that stimulus strength influence confidence differently as a function of the accuracy of the trial (Kepecs et al., 2008; Lak et al., 2014; Sanders et al., 2016). Navajas et al. (2017) provided a novel interpretation of these findings as they showed that computations corresponding to the probability to correct captured the well-established dissociation between error trials whereas computations capturing uncertainty in the evidence stream (the precision of the posterior belief) scaled positively with stimulus strength independent of first order accuracy. However, Navajas and colleagues also recorded individual differences in whether confidence judgments reflected probability correct or whether they reflected posterior precision. Additionally, the sample in the Kiani study consisted of 6 people. Thus there are two possible explanation of the Kiani effect: they had discovered evidence of distinct computations for sequential and simultaneous confidence judgments, or they happened to test participants who conceptualised confidence more as precision than the probability of being correct and misattributed this individual difference to their experimental design. Here I directly compared sequential and simultaneous confidence judgments from the same participants, and I have found no evidence that confidence scales positively with stimulus strength for error trials for either sequentially and simultaneously reported confidence, adding some plausibility to the individual differences account.

The within-participant design also allowed me to run a mediation analysis to test if the differential sensitivity to response time between sequential and simultaneous confidence judgments could explain the difference in metacognitive accuracy. I found that RT completely mediated the difference in the relationship between accuracy and confidence for the different

confidence timings. In other words, the difference in confidence between correct and incorrect trials was the same within people across conditions when the mediating effect of RT was controlled for. The sensitivity of confidence to response time differed between the sequential and simultaneous sessions, but did not vary between the sequential sessions as a function of interjudgment time, suggesting that there is a discrete step between simultaneous and sequential confidence judgments rather than a continuous shift based on increase processing time. These observations suggest that a discrete manipulation can lead to a graded effect (i.e. RT matters for confidence for both simultaneous and sequential responses, the effect is just stronger for the sequential ones). Interestingly this is not the first time this mix between a discrete manipulation and a graded effect has been reported in relation to confidence timings. Aitchinson and colleagues (2015) had a similar finding when they explored the computational underpinnings of confidence judgments. They tested which of three different computational substrates best captured the confidence judgments: the difference in stimulus strength between the chosen and unchosen option (the difference model), the stimulus strength of the chosen option only (the max model) or the probability of being correct as a function of the choice and the evidence for each option (the Bayesian model). They found that the Bayesian model fit the data best for sequential confidence judgments, but that a mix between the Bayesian model and the max model fitted the simultaneous response data best. In other words, in line with the current findings, their results indicate that the sequentially reported confidence judgments had a greater second-order accuracy than the simultaneously reported confidence judgments. Their findings may also be complementary with the current results. Kiani, Corthell and Shalden (2014) have argued that response time provides independent evidence of the probability that a choice is correct. Therefore, by combining my findings with the results of Aitchinson et al., one might speculate that sequential confidence, for most people, is an estimate of probability correct, and as such is partly dependent on response time. Simultaneous confidence judgments do not reflect probability correct, perhaps because the agent cannot compute the probability of being correct while the choice is being made. Instead simultaneous confidence judgments estimate $p(\text{correct})$ from a mixture of a partially completed Bayesian estimate and heuristic information. Because the sequential confidence report is completely based on a Bayesian computation, it is more accurate, and it also correlates more strongly to response time. Response time could easily be integrated into the Bayes-optimal computation of confidence because response time is a function of the area under the ramping curve in the decision phase. As such, a representation of response time can be captured from the same neurons that encode the evidence strength of the various options, provided that the decision threshold is fixed (Zylberberg, Barttfeld, & Sigman, 2012).

Readers might wonder about the discrepancy between my work and Aitchinson's work in that confidence for the sequential session in Aitchinson's model was based on a mixture between a model computing confidence from the evidence for the chosen option and a Bayesian model, whereas I computed stimulus strength as a difference. Indeed, other work has also suggested that the evidence for the unchosen option is irrelevant for confidence judgments (Maniscalco, Peters, & Lau, 2016; Zylberberg et al., 2012). In the case of Aitchinson's work the contrast between their findings and mine are mostly illusory, because simultaneous confidence in the Aitchinson study appeared to be a mixture of the max model and the Bayesian model, where the Bayesian model captures the influence of both negative and positive evidence. In fact, my results fit these findings well in that positive evidence predicted confidence more strongly than negative evidence in the current data, but the increased model fit relative to a difference model is not sufficient to justify the additional model complexity of adding an additional variable (see Appendix 4). Another difference between this work and experiments that have found that confidence is derived from positive evidence only is how the stimuli are presented. The experiments that only find an effect of positive evidence on confidence have presented the response options sequentially and under time constraints, as opposed to the free-viewing in the current experiment. This suggests that confidence judgments might rely more strongly on positive evidence in situations where memory-decay is a problem and response time is limited. In support of this theory, the second experiment reported by Zylberberg and colleagues (2012), which presents both options simultaneously and allows for free-viewing, shows a relationship between confidence and negative evidence, in line with my findings.

This study has extended previous work by showing that metacognitive accuracy is higher for sequentially reported confidence than simultaneously reported confidence within the same participants and that this finding can be explained by the additional influence of response time on the confidence judgments. This suggests that simultaneously and sequentially reported confidence differs, even within the same task. I explored whether this effect could be explained by differences in processing time, and while I could not rule out such an explanation, the evidence suggests that variations in additional total response time did not correspond to changes in metacognitive performance. I also replicated the role of eye behaviours in confidence and choice that I introduced in the previous chapter. These results suggest that the dynamics of visual attention play a similar role in perceptual and value-based decision making.

6. General Discussion

6.1. Summary

This chapter begins with a brief review of the main findings from the empirical chapters, followed by a discussion of their implications. With regard to the consequences of confidence, results from the value chapter imply that explicit confidence judgments serve a function to tag poor choices, which in turn enables agents to explore different options when they encounter a similar choice set in the future. This is contrasted from previous work on changes of mind that has studied changes occurring before a single trial is completed and as such might have more in common with error correction in perceptual decision making than confidence judgments and error monitoring. Ways to extend the current findings, by applying tasks that allow both for the objective evaluation of accuracy and for distinct choice sets that are recognisable by the participant, are discussed. The causes of confidence discuss two novel predictors of confidence, the summed value of the options presented in value-based choice and the amount of time the various options are fixated. This PhD suggests that eye behaviours relate to choice and confidence the same way in the perceptual domain and the value domain, suggesting that the dynamics of visual attention might have domain-general influences on decision making. Finally, the computational underpinnings of confidence are discussed in relation to the differences in metacognitive performance discovered between bilinguals and monolinguals, and between confidence judgments conducted simultaneously with or after choice.

6.2. Overview of Findings

The empirical chapters in this thesis have focused on three questions: (1) what factors influence confidence judgments, (2) how are confidence judgments computed, and (3) how do confidence judgments influence subsequent behaviours? The two experiments reported in Chapter 3 showed a previously unknown link between metacognitive accuracy and a trait: bilingualism. Specifically, the monolingual participants showed greater metacognitive accuracy than the bilingual participants in a visual discrimination task that controlled for first-order performance (and first-order response times in Experiment 2). The monolingual advantage in metacognitive accuracy was not caused by monolingual confidence judgments being more sensitive to response time or difficulty, nor could it be explained by a non-linear combination of the two. This opens up for the possibility that there are predictors of choice accuracy that inform confidence that have not yet been considered.

Experiments 3, 4 and 5 reported in Chapter 4 showed a potential candidate for such a novel predictor: Gaze shift frequency (GSF) – the number of times that participants shifted their visual attention between the available options – predicted between-trial fluctuations in confidence above and beyond choice difficulty and response time for both perceptual discriminations (Experiment 6) and value-based choices (Experiments 3 and 4). Additionally, GSF correlated with accuracy when response time and difficulty was controlled for (Experiment 6). Experiments 3 and 4 also showed that the summed value of the options positively predicted confidence whilst at the same time making participants *less likely* to pick the highest value option. This dissociation between confidence and performance will be further discussed in subsequent sections. Whereas Experiments 1 and 2 found a novel predictor of individual differences in metacognitive ability, Experiment 6 presented in Chapter 5 showed that changes in the timing of the confidence judgments relative to the choice alter participants' metacognitive sensitivity. Specifically, confidence judgments that followed first-order decisions were more strongly influenced by response time than confidence judgments that were made simultaneously with the choice. Response times were equally associated with accuracy for the sequential and simultaneous confidence judgments, hence sequential confidence judgments were more accurate than simultaneous confidence judgments because they captured the information in response time more effectively.

With regard to the consequences of confidence, Experiments 3 and 4 showed that confidence predicted changes of mind when the same choice was presented in subsequent trials. This extends the argument that confidence may serve as an error detection mechanism (Yeung & Summerfield, 2012, 2014) in the value domain when internal consistency of preferences serve as a proxy for accuracy. It also shows that confidence judgments may influence subsequent behaviours on a greater time scale than has been previously studied. Next I will explore this argument in a bit more detail.

6.3. The Consequences of Confidence

Much has been written about the function of confidence. On a neurological level it has been argued that representing the uncertainty of different information streams allows the brain to weight them based on their reliability, which results in better internal models (Beck et al., 2008; Ma, Beck, Latham, & Pouget, 2006). This basic idea can be extended to the level of explicit confidence of an agent, in that confidence judgments may enable an individual to weight information sources based on their reliability, in order to make more accurate decisions. Alternatively, confidence might help an individual determine when they should search for more

information prior to making a decision (Daw, O'Doherty, Dayan, Seymour, & Dolan, 2006). Explicit confidence representations may also support learning by enabling the agent to focus on the cues that are most predictive of the outcome (Meyniel & Dehaene, 2017). Furthermore, confidence may serve as a prediction error signal, enabling unsupervised learning in contexts where feedback is absent (Guggenmos et al., 2016). The ability to weight information sources may provide further benefits when we consider groups of individuals, as explicit confidence judgments could improve group decision making as the group assigns more weight to the accounts of the most confident members (Shea et al., 2014). However, this benefit may be mediated by the extent to which the group members have shared terminology for expressing their confidence estimates (Fusaroli et al., 2012).

One central function of confidence is allowing agents to change their minds. This aspect of confidence was first studied in the perceptual decision making literature in the context of error correction (Rabbitt, 1966) and was considered a separate area of enquiry until recently (Yeung & Summerfield, 2012, 2014). In perceptual decision making tasks it is common to distinguish between fast errors and slow errors (Pleskac & Busemeyer, 2010; Rabbitt, Cumming, & Vyas, 1978; Scheffers & Coles, 2000). Fast errors tend to occur under speed stress and may be resolved by motor-adjustment within the same trial. This form of online motor adjustment is called error correction in the perceptual decision making literature and does not seem to require conscious awareness (Nieuwenhuis, Ridderinkhof, Blom, Band, & Kok, 2001). In instances where fast errors are not corrected before an action is completed, it is often obvious to the agent that an error has been committed, perhaps because fast errors are the result of motor-interference rather than noisy evidence integration (Rabbitt et al., 1978; Scheffers & Coles, 2000). Slow errors, on the other hand, are not dependent on speed stress and are the result of noise in the stimulus or the decision making process. Slow errors can also be detected, but agents tend to be less certain as to whether they committed an error. Yeung and Summerfield (2012, 2014) showed the similarity between error monitoring of slow errors and confidence judgments and suggested they share the same mechanism. This assertion has since been supported empirically, as the EEG markers for error correction also predicted graded confidence judgments for correct trials (Boldt & Yeung, 2015).

Error monitoring for slow errors comes with several benefits apart from the opportunity to correct an immediate mistake (Resulaj et al., 2009). For example, it might allow agents to adjust their response threshold so that subsequent errors become less likely (Yeung & Summerfield, 2012). In the language of dynamic models of choice such as a drift diffusion models, the boundary separation might increase after an error, so that stronger evidence is required before an

action is initiated. This may explain the well-known phenomenon of post-error slowing (although the rareness of errors is almost certainly also a factor, see Notebaert et al., 2009). The relationship between error monitoring and boundary separation might be particularly valuable for complex tasks, where more than one step needs to be completed successfully to achieve a reward. Van den Berg and colleagues (2016) ran an experiment where participants completed pairs of visual discriminations. Participants were only rewarded if both discriminations were accurate. Using a computational model, the researchers found that participants adopted a laxer response criterion for the second trial if they reported low confidence in the first trial, so that they spent less time on trials that were unlikely to be rewarded. Together these results suggest that conscious error monitoring may allow the agent to adapt their boundary separation to maximise their chance of subsequent reward.

Previous research has illustrated many ways in which explicit confidence judgments might lead to better outcomes by influencing immediately subsequent decisions. Results from Experiments 3 and 4 in Chapter 4 extend this work by suggesting that confidence judgments may lead to better outcomes over time by causing changes of mind when the same options are encountered again. Perceptual decision making paradigms are poorly equipped to address this question because the stimuli in such tasks rarely have discrete identities, so even if the exact same choice were presented to a participant at multiple times during an experiment, it is unlikely that the participant would realise that they were encountering the same options. In Experiments 3 and 4 I addressed this issue in the domain of value-based choice by having hungry participants choose between common snack items. I found that participants who reported low confidence the first time they encountered an item set were more likely to change their mind the next time they encountered the same item set. Note that this definition of change of mind is quite different from the changes of mind that has typically been reported in the decision making literature, where people change their mind within the same trial as part of an ongoing decision process (Resulaj et al., 2009; Van Den Berg et al., 2016).

The fact that low confidence judgments predict changes of mind in future trials does not necessarily mean that confidence causes subsequent changes of mind. An alternative explanation follows. Low confidence trials are trials when the choice process was particularly noisy, because confidence acts as an error detection mechanism. This noise causes the worse of the two options being chosen at presentation one. When the same options are presented again the decision-process is less noisy and the best option is chosen, not as a result of the low confidence following the first choice but because the decision process for the second presentation will on average be less noisy because of regression to the mean. There are two reasons to be sceptical of

this account. First, I also found that highly metacognitive individuals became more internally consistent over time compared to their less metacognitive peers. It is unclear why this would happen if the confidence judgments did not influence subsequent decisions. Second, if we take the individual BDM-ratings to capture each person's true value representation of each item, low confidence predicts changes of mind regardless of whether the item with the highest BDM value was chosen originally, for both experiments. In other words, low confidence following the first encounter with a choice set leads to an increased chance of a change of mind, regardless of whether participants chose the "best option" or not.

As the previous section illustrates, the main challenge of applying the error-monitoring/error correction models from perceptual decision making to the value domain is that it is difficult to assess accuracy in the value context, because value is an inherently subjective quality. In this dissertation I used internal constancy as a proxy for accuracy, because inconsistent preference rankings allow an agent to be exploited and is therefore objectively suboptimal (Hájek, 2008; Von Neumann & Morgenstern, 2007). However, this approach is not perfect because while the algorithm used provided the most consistent item rankings for each participant, it did not guarantee that it was a unique solution (i.e. there might be many potential "best" rankings for each choice set; Pedings, Langville, & Yamamoto, 2012). Additionally, it is possible that a participant genuinely changed their preferences during the course of the experiment; it is impossible for an outside observer to differentiate between an "error" where a participant picked a less preferred option by mistake from a "change of mind" where a participant genuinely changed their preference.

Fully investigating the relationship between confidence, long-term changes of mind and error monitoring would require two things. First, each option presented needs to be identifiable to the participant so that they are aware that choice sets repeat. Most perceptual discrimination tasks fail this criterion. Second, each choice needs to have an objectively correct option so that the accuracy of the initial choice can be controlled for when assessing the role of confidence on subsequent changes of mind. The value tasks presented in this dissertation fail this criterion because value is inherently subjective. One task that would fulfil both these criteria would be a memory task where participants learn a set of word associations, and then have to pick the target word linked to a cue from a set of options. After each choice participants would indicate their confidence in their accuracy, and after a set of trials the same choice would recur. Such an experiment would extend our knowledge of the role of confidence in long-term-changes of mind outside of the value domain. Additionally it would provide a stronger causal test of the role of confidence in changes of mind, as the influence of confidence could be evaluated when

controlling for the difficulty of the item, the response time of the participant and the accuracy of the original choice.

6.4. The Causes of Confidence

It has long been known that confidence is associated both with response time and stimulus strength (Baranski & Petrusic, 1994, 1998; De Martino et al., 2013; Festinger, 1943; Kiani & Shadlen, 2009; Ratcliff & Starns, 2009; Vickers & Packer, 1982) and recently Kiani, Corthell and Shadlen (2014) demonstrated that response time has a causal influence on confidence judgments. The relationship between confidence, response time and trial difficulty is often explained in relation to sequential sampling models, as a low drift rate, which is associated with difficult trials would result in both slower responding and more errors (Pleskac & Busemeyer, 2010; Vickers & Packer, 1982). Recent work in the value domain suggests that drift rate is influenced by visual attention (Ian Krajbich et al., 2010; Ian Krajbich & Rangel, 2011; Lim et al., 2011). The original work from Rangel's lab suggested that the time spent looking at the item interacted with the item value to determine drift rate, but subsequent work suggests that the effect is in fact additive in that the fixated item gets a small constant boost, independent of its value (Cavanagh et al., 2014). I have replicated this additive effect of fixation time in both the value domain (Experiment 3) and in relation to perceptual choice (Experiment 6). The replication in the value domain is important because the items used for the value-based experiment was much closer to Krajbich (who also used snack items) than Cavanagh (who used bandits whose values had been taught in the same experiment), so the results suggest that visual attention influences drift rate in an additive fashion even when the participants have extensive experience with the items. The results from the perceptual task suggest that fixation time additively influences drift rate for visual discriminations as well, suggesting that this effect might be domain general. Given the effect of visual attention on evidence accumulation, I wanted to examine how it relates to confidence.

I did not find a reliable relationship between fixation time and confidence in value-based or perceptual choice (Experiments 3, 4 and 5). However, the number of times that a participant shifted their gaze between the options (Gaze Shift Frequency, GSF) negatively predicted confidence. The influence of GSF on confidence was reliable even when response time and trial difficulty was controlled for, in both perceptual judgments and for value-based choice. GSF is a novel measure of uncertainty that is both conceptually and empirically independent from differences in total fixation time. It is currently unclear if GSF is associated with confidence because an internal representation of GSF feeds into the confidence judgment or whether participants simply move their eyes more when they are more uncertain, and this uncertainty is

also reflected in their confidence judgment. In other words, it is currently not clear if GSF causally influences confidence, or whether both GSF and confidence are driven by an internal uncertainty representation. The easiest way to establish the causal relationship between GSF and confidence would be to run an experiment where GSF is independently manipulated while keeping other important variables (notably difficulty and response time) constant. If externally manipulated GSF still predicts subsequent confidence judgments it would suggest that it has a causal influence.

Another interesting question is whether GSF is specific to the visual modality or whether it is a more general mechanism for how humans sample evidence from a set of options. This question can be further subdivided into a motor component and a perception component. Specifically, is there something about eye movements that is associated with confidence, or alternatively, is there something about visual evidence sampling that is associated with confidence? Both these questions could be investigated in a single perceptual two-alternative-forced-choice task with a 2-by-2 design. One dimension would be the modality of the perceptual judgment (for example sight and sound). Each option would be sampled sequentially so that participants can shift between the modalities with a motor action. The second dimension would be the type of motor action; evidence sampling would either be controlled by ocular movement or some other motor action (e.g. pressing a button). This would result in 4 conditions:

- (1) Visual discrimination trials where the options are sampled by ocular motor action (similar to the perceptual discrimination task in Experiment 6)
- (2) Visual discrimination trials where changing between the options does not require ocular movements (e.g. one option is shown at the middle of the screen, and this option is replaced by a button press)
- (3) Auditory discrimination trials where options are sampled by an ocular action (for example a pitch discrimination where the participant has to determine which of two tones are the lowest and the tone heard is determined by where participants look on the screen)
- (4) Auditory discrimination trials where changing between the options does not require ocular movements (for example a pitch discrimination where the tone heard is determined by button press)

If the number of times participants shift between sampling each option predicts confidence across all conditions, that would suggest that the shift rate between the options matter, regardless

of modality. If only the eye movement conditions (1 & 3) were associated with confidence, that would suggest that confidence is informed by some form of motor-readout that is specific to the visual system. If only the visual perception conditions (1 & 2) found an association between shifting and confidence, that would imply that there is a special relationship between confidence and visual attention in humans. Finally, if the relationship between confidence and shifting only appeared for the first condition, that would suggest that this relationship requires both visual input and ocular motor output.

If GSF is proven to be robust and domain-general it could be a useful measure of uncertainty in non-human animals. This would be a valuable methodological contribution as previous work in animals has measured uncertainty as post-decision wagering (Kepecs & Mainen, 2012) or waiting time (Lak et al., 2014). The downside of post-decision wagering is that it behaviourally combines the choices animals make with their confidence judgments and therefore conflates uncertainty estimates with economic preferences such as loss aversion (Fleming & Dolan, 2010).

Additionally, post-decision wagering is a binary uncertainty-driven behaviour and as such is poorly suited to evaluate parametric models of confidence. Lak and colleagues (2014) invented a parametric measure of confidence in non-human animals. After a successful perceptual discrimination, the animal was rewarded after a random time interval, during this interval the animal could at any time start a new trial. Because the waiting time is associated with an opportunity cost, animals should only chose to wait on trials when they are certain they are correct and thus certain they will receive a reward, and the amount of time they are willing to wait should scale with their confidence. Unfortunately, waiting time is a noisy measure of confidence for most correct trials, as it is cut short whenever the reward appears. Consequently, Lak and colleagues withheld rewards on a small proportion of the correct trials, and used these uninterrupted waiting times as a measure of parametric confidence, meaning that only a small amount of the total choice trials came with a confidence estimate. GSF would address all of these problems as it would be a parametric measure of uncertainty that would be available for all trials (correct and incorrect) without being conflated with economic considerations. Consequently it would be well suited as a behavioural marker of uncertainty in the search for the neural substrates of decision confidence.

GSF is not the only new predictor of confidence presented in this dissertation. In Experiments 3 and 4 I showed that the summed value of the presented options positively predicted confidence when other predictors such as the value-difference between the options, the response time and GSF was accounted for. Summed-value is interesting because it is associated with both increased confidence and with making choices less sensitive to the value difference between the options.

To the best of my knowledge, this is the first time a variable in value-based choice has been shown to be associated with both increased confidence and decreased choice quality. There are two explanations that might account for the effect of summed value, one relating to its role in confidence, and the other relating to its role in choice quality. First, summed value might lead to higher confidence because the value signal might somehow contaminate the confidence signal. Neuroscientific work lends some credence to this idea as vmPFC – which has been implicated in value computations – also appears to be involved in confidence judgments (De Martino et al., 2013; Lebreton et al., 2015). Second, total value might decrease the sensitivity of the value difference in predicting choice because value difference is actually encoded as a ratio of the total value (e.g. a value difference of £0.5 is relatively smaller when the options are valued at £2.5 and £3.0 than when they are valued at £0.5 and £0). This fits the notion of divisive normalisation that has been observed in some neural codes (Carandini & Heeger, 2012; Louie et al., 2013; Soltani et al., 2012). Neural evidence could help test both these hypothesis. In order to explore these effect one could design choice sets with much more pronounced differences in total value than the ones presented here. Choices could then be explored both in terms of behavioural outcomes and with regards to BOLD response in the vmPFC.

6.5. The Computation of Confidence

This dissertation has not only introduced novel predictors of confidence, it also contains a novel moderator of existing predictors. Specifically, Experiment 6 in Chapter 5 showed that the influence of response time on confidence is moderated by whether the confidence judgment is reported simultaneously with the choice or after the choice. Previous work has explored how the timing of the confidence judgments influences the underlying computation (Aitchison et al., 2015; Kiani et al., 2014; Siedlecka et al., 2016), but they have all relied on between-participant comparisons. Because of reliable individual differences in metacognitive abilities (Ais et al., 2016) and processing styles (Navajas et al., 2017) any between-participant comparison is problematic. To the best of my knowledge, Experiment 6 is the first study to directly assess the influence of the timing of confidence judgments within the same participants. This allowed me to account for any variation due to individual differences and therefore made the results more reliable, and their interpretation more straightforward compared to previous work. I found that sequential confidence judgments were more accurate than simultaneous confidence judgments. These results are in line with the findings of Siedlecka and colleagues (2016) who compared the accuracy of confidence judgments conducted before and after a first order decision and who found that retrospective confidence judgments were more accurate. One potential explanation

for this difference in performance is put forward by Aitchison and colleagues (2015) who found that sequential confidence judgments mostly reflect the probability correct ($p(\text{correct})$) whereas simultaneous confidence judgments are also influenced by simpler heuristics (in their case the sensory magnitude of the chosen option). I found that the influence of response time on confidence was weaker for the sequential confidence judgments and that this difference in sensitivity to response time fully explained the difference in metacognitive efficiency. It is important to note that these explanations are not mutually exclusive as response time is predictive of $p(\text{correct})$ (Kiani et al., 2014). It is possible that it is challenging for the brain to compute $p(\text{correct})$ before an action has been taken (Fleming & Daw, 2017), and that it therefore relies on heuristics until a decision has been made. Once the first-order decision has been made the probability of being correct can be computed, and reaction time can be taken into account in this computation.

One way to evaluate the idea that simultaneous and sequential confidence judgments reflect different computational quantities is to look at the relationship between confidence and stimulus strength for correct trials and error trials. Navajas and colleagues (2017) showed that confidence judgments based on $p(\text{correct})$ show a positive association between confidence magnitude and stimulus strength for correct trials but a negative relationship for errors. However, confidence judgments that reflect the precision of the posterior evidence distribution (a heuristic approximating $p(\text{correct})$) show a positive association between confidence and stimulus strength for correct trials and for error trials. Kiani and Shadlen (2014) noted that stimulus strength positively predicted confidence both for the correct trials and the error trials for the participants who made simultaneous confidence judgments and choices, in conflict with much previous research on sequential confidence judgments (Kepecs et al., 2008; Lak et al., 2014; Sanders et al., 2016). In Experiment 6 I found a positive relationship between confidence and stimulus strength for correct trials but a negative association between confidence and stimulus strength for error trials, with no evidence that the timing of the confidence judgments moderated this pattern. Navajas and colleagues (2017) offer a simple explanation of these discrepancies, they found that confidence reflected some mixture of probability correct and the posterior precision for about 47% of their participants, whereas it just reflected probability correct for an additional 43%, they also found that the tendency to base confidence on specific computational quantities were stable across time. Because Kiani and Shadlen only tested six participants in their simultaneous confidence judgment task, it is possible that these six participants happened to base their confidence judgments primarily on the posterior precision. Thus it is possible that the small sample size, together with the between-participant design made the researchers misattribute an

effect that was caused by individual differences to their research manipulation. Such misattributions could have been avoided if they had compared the performance on both tasks within the same participants, further highlighting the importance of within subject comparisons when evaluating the influence of response timings on confidence judgments.

My results also include a difference in metacognitive performance that is not due to an experimental manipulation, but rather due to a group difference. Specifically I found that bilinguals were less metacognitively accurate than monolinguals in a visual discrimination task, despite comparable first-order performance between both groups (Experiments 1 and 2). This difference could not be explained by confidence being more sensitive to difficulty or response time in the monolingual group. In fact, these two variables seemed to influence confidence the same way in both groups. Future research is required to explain this difference. One potential explanation would be difference in sensitivity to GSF. The plausibility of this idea depends on the magnitude of the GSF effect. In Experiment 6, which used a similar design to the bilingual experiments, GSF was the strongest predictor of confidence. However, stimulus presentation was gaze contingent in Experiment 6 and it is possible that this manipulation inflated the size of the effect. Fleming and Lau (2017) has suggested that the choice itself can provide information to confidence judgments if the internal evidence streams causing confidence and choice are at least partially uncorrelated. This opens for the possibility that the metacognitive insight is better among monolinguals because their evidence streams for confidence and choice are more independent than their bilingual peers. This could be evaluated by having monolinguals and bilinguals completing two versions of a perceptual discrimination task where they either have to provide confidence judgments before or after they report their choice. If the bilingual disadvantage disappears for trials where confidence judgment precede choice but persist for trials where confidence judgments follow choice this would suggest that monolingual confidence judgments are more accurate because their choice and confidence evidence streams are more decoupled.

Additionally, because bilinguals reported systematically higher confidence ratings in Experiments 1 and 2 it would be interesting to test which group is better calibrated. This could be tested in any follow-up experiment just by presenting the confidence scale explicitly in terms of probability correct. Finally, it is possible that I failed to discover the cause of the performance difference between monolinguals and bilinguals because dot difference is a coarse measure of stimulus strength. The difficulty of the dot discrimination task also depends on the distribution of the dots: trials are harder when the dots are clustered or overlap. Navajas and colleagues developed an elegant experimental design that provides stricter control over a set of parameters

that influence difficulty. They present participants with a series of Gabor patches drawn from a uniform distribution with a mean tilt either to the left or to the right. Applying such a design to bilingual and monolingual samples would allow them to be compared in terms of their sensitivity to the magnitude of the evidence (difference in degrees between the left-tilting samples and the right-tilting samples), evidence variance (the width of the uniform distributions) and number of evidence samples. Such a design has the additional benefit that it would be easy to compute the posterior distribution of the tilt at the time of choice (assuming a flat prior at the start of each trial) and would therefore allow for a direct test of the extent to which confidence is a function of posterior precision relative to probability correct for monolingual and bilingual participants. It is worth pointing out that if such a difference exists it is likely to be small because it is not captured by the relationship between confidence and stimulus strength for correct trials and error trials in either of the experiments in Chapter 3 (see Figure 3.9.).

Some models of confidence view it as the result of the continuation of the same evidence accumulation process that cause first-order decision (Kiani et al., 2014; Pleskac & Busemeyer, 2010; Van Den Berg et al., 2016). I mentioned in the introduction that these theories cannot provide a complete picture of confidence judgments because they cannot account for differences in metacognitive accuracy for trials with the same response times (Vlassova et al., 2014). They also struggle to account for effects where a variable influences first-order and second order accuracy in opposite directions, as was the case for summed value, reported in Chapter 4. Single stream theories also have trouble with order effects, i.e. how confidence judgments alter based on whether they happen before or after a first-order decision (Fleming & Daw, 2017; Siedlecka et al., 2016). Experiment 6 in this thesis records another such order effect, but fails to conclusively rule out that this order effect is caused by differences in processing time or cognitive load. Together, these findings leave us with a dilemma; some results suggest that a single evidence stream best explains confidence and choice, whereas other results are incongruent with such models. One potential solution comes from the research on error detection: in their review paper Yeung and Summerfield (2012) point out that there are two error-sensitive signals in the brain, one that emphasise immediate error correction and one that emphasise error monitoring. In real world contexts where tasks are less clearly delineated than the lab it is possible that the error-correcting code does not just correct immediate errors but supports online motor adjustments for activities like walking or riding a bike, that require constant calibration between incoming sensory data and motor output. On the other side of the spectrum we have an error-monitoring/confidence signal that operates on discrete actions and judge their probability of success in the service of learning and strategic decision making (e.g.

how long to wait for a reward, whether to adjust response boundaries to be more careful and accurate, whether to keep exploiting a current resources or explore the environment). It would be adaptive for this online motor code to be a direct extension of the evidence accumulation process that controls the initial action because it is a natural extension of that process (in fact, it might be exactly the same process, arbitrarily delimited by the tasks we use in the lab). On the other hand, the slower system that explicitly monitors the probability of success likely transforms the sensory evidence signal into a probability estimate. If current research on confidence conflates these two signals, it would explain many of the seemingly conflicting findings we see today. These two signals might be disentangled with neuroimaging techniques as suggested by Yeung and Summerfield since error correction and error monitoring are associated with different EEG signals (ERN and Pe respectively). Using EEG to capture the input of these two components into confidence judgments in various experimental setups would be a useful research project. Only once we have a better sense of what neural and computational quantities underlie self-reported confidence in different contexts can we gain a complete understanding of the function of confidence.

6.6. Conclusion

Uncertainty is an inherent property of the world and of the neurological system we use to model it. Therefore understanding uncertainty is central for successfully negotiating our existence. This thesis has explored our ability to consciously quantify uncertainty as confidence, and what benefits we derive from doing so. It has uncovered a number of novel predictors of confidence: internal predictors like eye behaviours, situational predictors like the total value of the options under consideration and interpersonal, and trait-like predictors such as the number of languages one speaks. It has made a novel suggestion for the function of confidence, namely that it tags poor decisions so that when similar situations arise in the future different options may be chosen. Finally, it has added to our understanding of the computational underpinnings of confidence by demonstrating that confidence judgments capture different information based on how they are timed in relation to choice. Together these findings support an emerging consensus that self-reported confidence judgments correspond to different computational quantities as a function of task structure as well as individual differences. Future research must continue to map what computational and neural quantities underlie different forms of confidence judgments. Only when that mapping is clearer can a complete model of the pragmatics of confidence be attempted.

References

- Abutalebi, J., & Green, D. (2007). Bilingual language production: The neurocognition of language representation and control. *Journal of Neurolinguistics*, 20(3), 242–275.
- Ackerman, R. (2014). The diminishing criterion model for metacognitive regulation of time investment. *Journal of Experimental Psychology: General*, 143(3), 1349.
- Ais, J., Zylberberg, A., Barttfeld, P., & Sigman, M. (2016). Individual consistency in the accuracy and distribution of confidence judgments. *Cognition*, 146, 377–386.
- Aitchison, L., Bang, D., Bahrami, B., & Latham, P. E. (2015). Doubly Bayesian analysis of confidence in perceptual decision-making. *PLoS Comput Biol*, 11(10), e1004519.
- Antón, E., Duñabeitia, J. A., Estévez, A., Hernández, J. A., Castillo, A., Fuentes, L. J., ... Carreiras, M. (2014). Is there a bilingual advantage in the ANT task? Evidence from children. *Frontiers in Psychology*, 5, 398.
- Bach, D. R., & Dolan, R. J. (2012). Knowing how much you don't know: a neural organization of uncertainty estimates. *Nature Reviews Neuroscience*, 13(8), 572–586.
- Badre, D., Doll, B. B., Long, N. M., & Frank, M. J. (2012). Rostrolateral prefrontal cortex and individual differences in uncertainty-driven exploration. *Neuron*, 73(3), 595–607.
- Bahrami, B., Olsen, K., Bang, D., Roepstorff, A., Rees, G., & Frith, C. (2012). What failure in collective decision-making tells us about metacognition. *Phil. Trans. R. Soc. B*, 367(1594), 1350–1365.
- Baird, B., Cieslak, M., Smallwood, J., Grafton, S. T., & Schooler, J. W. (2015). Regional white matter variation associated with domain-specific metacognitive accuracy. *Journal of Cognitive Neuroscience*, 27(3), 440–452.
- Baird, B., Smallwood, J., Gorgolewski, K. J., & Margulies, D. S. (2013). Medial and lateral networks in anterior prefrontal cortex support metacognitive ability for memory and perception. *Journal of Neuroscience*, 33(42), 16657–16665.
- Bang, D., Fusaroli, R., Tylén, K., Olsen, K., Latham, P. E., Lau, J. Y. F., ... Bahrami, B. (2014). Does interaction matter? Testing whether a confidence heuristic can replace interaction in collective decision-making. *Consciousness and Cognition*, 26, 13–23.
- Baranski, J. V., & Petrusic, W. M. (1994). The calibration and resolution of confidence in perceptual judgments. *Perception & Psychophysics*, 55(4), 412–428.
- Baranski, J. V., & Petrusic, W. M. (1998). Probing the locus of confidence judgments: experiments on the time to determine confidence. *Journal of Experimental Psychology: Human Perception and Performance*, 24(3), 929–945.
- Baranski, J. V., & Petrusic, W. M. (2001). Testing architectures of the decision–confidence relation. *Canadian Journal of Experimental Psychology/Revue Canadienne de Psychologie Expérimentale*, 55(3), 195–206.
- Barrett, A. B., Dienes, Z., & Seth, A. K. (2013). Measures of metacognition on signal-detection theoretic models. *Psychological Methods*, 18(4), 535–552.
- Barron, H. C., Garvert, M. M., & Behrens, T. E. J. (2015). Reassessing VMPFC: full of

- confidence? *Nature Neuroscience*, 18(8), 1064–1066.
- Basten, U., Biele, G., Heekeren, H. R., & Fiebach, C. J. (2010). How the brain integrates costs and benefits during decision making. *Proceedings of the National Academy of Sciences*, 107(50), 21767–21772.
- Bates, D., Maechler, M., Bolker, B., & Walker, S. (2014). lme4: Linear mixed-effects models using Eigen and S4. *Journal of Statistical Software*. Retrieved from <http://arxiv.org/abs/1406.5823>
- Beck, J. M., Ma, W. J., Kiani, R., Hanks, T., Churchland, A. K., Roitman, J., ... Pouget, A. (2008). Probabilistic population codes for Bayesian decision making. *Neuron*, 60(6), 1142–1152.
- Becker, G. M., DeGroot, M. H., & Marschak, J. (1964). Measuring utility by a single-response sequential method. *Behavioral Science*, 9(3), 226–232.
- Begg, I., Duft, S., Lalonde, P., Melnick, R., & Sanvito, J. (1989). Memory predictions are based on ease of processing. *Journal of Memory and Language*, 28(5), 610–632.
- Benjamin, A. S., Bjork, R. A., & Schwartz, B. L. (1998). The mismeasure of memory: when retrieval fluency is misleading as a metamnemonic index. *Journal of Experimental Psychology: General*, 127(1), 55.
- Bialystok, E. (2001). *Bilingualism in development: Language, literacy, and cognition*. Cambridge University Press.
- Bialystok, E., Craik, F. I. M., Klein, R., & Viswanathan, M. (2004). Bilingualism, aging, and cognitive control: evidence from the Simon task. *Psychology and Aging*, 19(2), 290–303.
- Bialystok, E., & Martin, M. M. (2004). Attention and inhibition in bilingual children: Evidence from the dimensional change card sort task. *Developmental Science*, 7(3), 325–339.
- Bialystok, E., & Viswanathan, M. (2009). Components of executive control with advantages for bilingual children in two cultures. *Cognition*, 112(3), 494–500.
- Bogacz, R., Brown, E., Moehlis, J., Holmes, P., & Cohen, J. D. (2006). The physics of optimal decision making: a formal analysis of models of performance in two-alternative forced-choice tasks. *Psychological Review*, 113(4), 700–765.
- Boldt, A., Blundell, C., & De Martino, B. (2017). Confidence modulates exploration and exploitation in value-based learning. *bioRxiv*, 236026.
- Boldt, A., De Gardelle, V., & Yeung, N. (2017). The impact of evidence reliability on sensitivity and bias in decision confidence. *Journal of experimental psychology: human perception and performance*, 43(8), 1520.
- Boldt, A., & Yeung, N. (2015). Shared neural markers of decision confidence and error detection. *Journal of Neuroscience*, 35(8), 3478–3484.
- Bolker, B. (2014, October). How trustworthy are the confidence intervals for lmer objects through effects package? Retrieved from <http://stats.stackexchange.com/questions/117641/how-trustworthy-are-the-confidence-intervals-for-lmer-objects-through-effects-pa>
- Boorman, E. D., Behrens, T. E. J., Woolrich, M. W., & Rushworth, M. F. S. (2009). How green is the grass on the other side? Frontopolar cortex and the evidence in favor of alternative courses of action. *Neuron*, 62(5), 733–743.
- Botvinick, M. M., Cohen, J. D., & Carter, C. S. (2004). Conflict monitoring and anterior cingulate

- cortex: an update. *Trends in Cognitive Sciences*, 8(12), 539–546.
- Bronfman, Z. Z., Brezis, N., Moran, R., Tsetsos, K., Donner, T., & Usher, M. (2015). Decisions reduce sensitivity to subsequent information. *Proceedings of the Royal Society of London Series Biological Sciences*, 282, 20150228.
- Camerer, C. F., & Ho, T.-H. (1994). Violations of the betweenness axiom and nonlinearity in probability. *Journal of Risk and Uncertainty*, 8(2), 167–196.
- Carandini, M., & Heeger, D. J. (2012). Normalization as a canonical neural computation. *Nature Reviews Neuroscience*, 13(1), 51–62.
- Carlson, S. M., & Meltzoff, A. N. (2008). Bilingual experience and executive functioning in young children. *Developmental Science*, 11(2), 282–298.
- Carroll, M., Nelson, T. O., & Kirwan, A. (1997). Tradeoff of semantic relatedness and degree of overlearning: Differential effects on metamemory and on long-term retention. *Acta Psychologica*, 95(3), 239–253.
- Cavanagh, J. F., Wiecki, T. V., Kochar, A., & Frank, M. J. (2014). Eye tracking and pupillometry are indicators of dissociable latent decision processes. *Journal of Experimental Psychology. General*, 143(4), 1476–88.
- Chandler, C. C. (1994). Studying related pictures can reduce accuracy, but increase confidence, in a modified recognition test. *Memory & cognition*, 22(3), 273–280.
- Chib, V. S., Rangel, A., Shimojo, S., & O'Doherty, J. P. (2009). Evidence for a common representation of decision values for dissimilar goods in human ventromedial prefrontal cortex. *Journal of Neuroscience*, 29(39), 12315–12320.
- Cohen, J. D., McClure, S. M., & Angela, J. Y. (2007). Should I stay or should I go? How the human brain manages the trade-off between exploitation and exploration. *Philosophical Transactions of the Royal Society of London B: Biological Sciences*, 362(1481), 933–942.
- Costa, A., Hernández, M., Costa-Faidella, J., & Sebastián-Gallés, N. (2009). On the bilingual advantage in conflict processing: Now you see it, now you don't. *Cognition*, 113(2), 135–149.
- Costa, A., Hernández, M., & Sebastián-Gallés, N. (2008). Bilingualism aids conflict resolution: Evidence from the ANT task. *Cognition*, 106(1), 59–86.
- Couchman, J. J., Coutinho, M. V., Beran, M. J., & Smith, J. D. (2010). Beyond stimulus cues and reinforcement signals: a new approach to animal metacognition. *Journal of Comparative psychology*, 124(4), 356.
- David, A. S., Bedford, N., Wiffen, B., & Gilleen, J. (2012). Failures of metacognition and lack of insight in neuropsychiatric disorders. *Phil. Trans. R. Soc. B*, 367(1594), 1379–1390.
- Daw, N. D., O'doherty, J. P., Dayan, P., Seymour, B., & Dolan, R. J. (2006). Cortical substrates for exploratory decisions in humans. *Nature*, 441(7095), 876–879.
- de Bruin, A., Treccani, B., & Della Sala, S. (2015). Cognitive advantage in bilingualism: an example of publication bias? *Psychological Science*, 26(1), 99–107.
- De Gardelle, V., & Mamassian, P. (2014). Does confidence use a common currency across two visual tasks? *Psychological Science*, 25(6), 1286–1288.
- De Martino, B., Fleming, S. M., Garrett, N., & Dolan, R. J. (2013). Confidence in value-based choice. *Nature Neuroscience*, 16(1), 105–110.

- Dehaene, S., & Sigman, M. (2012). From a single decision to a multi-step algorithm. *Current Opinion in Neurobiology*, 22(6), 937–945.
- Del Missier, F., Mäntylä, T., & Bruine de Bruin, W. (2010). Executive functions in decision making: An individual differences approach. *Thinking & Reasoning*, 16(2), 69–97.
- Dijkstra, T., De Bruijn, E., Schriefers, H., & Ten Brinke, S. (2000). More on interlingual homograph recognition: Language intermixing versus explicitness of instruction. *Bilingualism: Language and Cognition*, 3(1), 69–78.
- Evans, S., & Azzopardi, P. (2007). Evaluation of a 'bias-free' measure of awareness. *Spatial Vision*, 20(1), 61–77.
- Fernandez-Duque, D., Baird, J. A., & Posner, M. I. (2000). Executive attention and metacognitive regulation. *Consciousness and Cognition*, 9(2), 288–307.
- Festinger, L. (1943). Studies in decision: I. Decision-time, relative frequency of judgment and subjective confidence as related to physical stimulus difference. *Journal of Experimental Psychology*, 32(4), 291–306.
- Filippi, R., Leech, R., Thomas, M. S. C., Green, D. W., & Dick, F. (2012). A bilingual advantage in controlling language interference during sentence comprehension. *Bilingualism: Language and Cognition*, 15(4), 858–872.
- Filippi, R., Morris, J., Richardson, F. M., Bright, P., Thomas, M. S. C., Karmiloff-Smith, A., & Marian, V. (2015). Bilingual children show an advantage in controlling verbal interference during spoken language comprehension. *Bilingualism: Language and Cognition*, 18(3), 490–501.
- FitzGerald, T. H. B., Seymour, B., & Dolan, R. J. (2009). The role of human orbitofrontal cortex in value comparison for incommensurable objects. *Journal of Neuroscience*, 29(26), 8388–8395.
- Flavell, J. H. (1979). Metacognition and cognitive monitoring: A new area of cognitive–developmental inquiry. *American Psychologist*, 34(10), 906–911.
- Fleming, S. M. (2017). HMeta-d: hierarchical Bayesian estimation of metacognitive efficiency from confidence ratings. *Neuroscience of Consciousness*, 3(1).
- Fleming, S. M., & Daw, N. D. (2017). Self-evaluation of decision-making: A general Bayesian framework for metacognitive computation. *Psychological Review*, 124(1), 91–114.
- Fleming, S. M., & Dolan, R. J. (2010). Effects of loss aversion on post-decision wagering: implications for measures of awareness. *Consciousness and Cognition*, 19(1), 352–363.
- Fleming, S. M., & Dolan, R. J. (2012). The neural basis of metacognitive ability. *Philosophical Transactions of the Royal Society of London B: Biological Sciences*, 367(1594), 1338–1349.
- Fleming, S. M., Dolan, R. J., & Frith, C. D. (2012). Metacognition: computation, biology and function. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 367(1594), 1280–1286.
- Fleming, S. M., Huijgen, J., & Dolan, R. J. (2012). Prefrontal contributions to metacognition in perceptual decision making. *The Journal of Neuroscience*, 32(18), 6117–6125.
- Fleming, S. M., & Lau, H. C. (2014). How to measure metacognition. *Frontiers in Human Neuroscience*, 8, 443.
- Fleming, S. M., Ryu, J., Golfinos, J. G., & Blackmon, K. E. (2014). Domain-specific impairment in metacognitive accuracy following anterior prefrontal lesions. *Brain*, 137(10), 2811–2822.

- Fleming, S. M., Weil, R. S., Nagy, Z., Dolan, R. J., & Rees, G. (2010). Relating introspective accuracy to individual differences in brain structure. *Science*, 329(5998), 1541–1543.
- Folke, T., Jacobsen, C., Fleming, S. M., & De Martino, B. (2016). Explicit representation of confidence informs future value-based decisions. *Nature Human Behaviour*, 1, 2.
- Folke, T., Ouzia, J., Bright, P., De Martino, B., & Filippi, R. (2016). A bilingual disadvantage in metacognitive processing. *Cognition*, 150, 119–132.
- Foot, A. L., & Crystal, J. D. (2007). Metacognition in the rat. *Current Biology*, 17(6), 551–555.
- Forstmann, B. U., Ratcliff, R., & Wagenmakers, E. J. (2016). Sequential sampling models in cognitive neuroscience: Advantages, applications, and extensions. *Annual review of psychology*, 67, 641–666.
- Friedman, J., Hastie, T., & Tibshirani, R. (2001). *The elements of statistical learning*. Springer series in statistics Springer, Berlin.
- Fujita, K. (2009). Metamemory in tufted capuchin monkeys (*Cebus apella*). *Animal cognition*, 12(4), 575.
- Fusaroli, R., Bahrami, B., Olsen, K., Roepstorff, A., Rees, G., Frith, C., & Tylén, K. (2012). Coming to terms: quantifying the benefits of linguistic coordination. *Psychological Science*, 23(8), 931–939.
- Galvin, S. J., Podd, J. V., Draga, V., & Whitmore, J. (2003). Type 2 tasks in the theory of signal detectability: Discrimination between correct and incorrect decisions. *Psychonomic Bulletin & Review*, 10(4), 843–876.
- García, G. E., Jiménez, R. T., & Pearson, P. D. (1998). Metacognition, childhood bilingualism, and reading. *Metacognition in Educational Theory and Practice*, 193–219.
- Gathercole, V. C. M., Thomas, E. M., Kennedy, I., Prys, C., Young, N., Viñas-Guash, N., ... Jones, L. (2014). Does language dominance affect cognitive performance in bilinguals? Lifespan evidence from preschoolers through older adults on card sorting, Simon, and metalinguistic tasks. *Frontiers in Psychology*, 5, 11.
- Gelman, A., & Hill, J. (2006). *Data analysis using regression and multilevel/hierarchical models*. Cambridge University Press.
- Gherman, S., & Philiastides, M. G. (2015). Neural representations of confidence emerge from the process of decision formation during perceptual choices. *Neuroimage*, 106, 134–143.
- Gilbert, S. J. (2015). Strategic use of reminders: Influence of both domain-general and task-specific metacognitive confidence, independent of objective memory ability. *Consciousness and Cognition*, 33, 245–260.
- Gold, J. I., & Shadlen, M. N. (2007). The neural basis of decision making. *Annu. Rev. Neurosci.*, 30, 535–574.
- Graziano, M., & Sigman, M. (2009). The spatial and temporal construction of confidence in the visual scene. *PLoS One*, 4(3), e4909.
- Green, D., & Swets, J. (1996). *Signal detection theory and psychophysics*. New York: Wiley.
- Green, D. W. (1986). Control, activation, and resource: A framework and a model for the control of speech in bilinguals. *Brain and Language*, 27(2), 210–223.

- Green, D. W. (1998). Mental control of the bilingual lexico-semantic system. *Bilingualism: Language and Cognition*, 1(2), 67–81.
- Grimaldi, P., Lau, H., & Basso, M. A. (2015). There are things that we know that we know, and there are things that we do not know we do not know: Confidence in decision-making. *Neuroscience & Biobehavioral Reviews*, 55, 88–97.
- Guggenmos, M., Wilbertz, G., Hebart, M. N., & Sterzer, P. (2016). Mesolimbic confidence signals guide perceptual learning in the absence of external feedback. *eLife*, 5, e13388.
- Hájek, A. (2008). Dutch book arguments. In P. Anand, P. Pattanaik, & C. Puppe (Eds.), *The Oxford Handbook of Rational and Social Choice*, (pp. 173–195).
- Halekoh, U., & Højsgaard, S. (2014). A Kenward-Roger approximation and parametric bootstrap methods for tests in linear mixed models—the R package pbkrtest. *Journal of Statistical Software*, 59(9), 1–32.
- Hampton, R. R. (2001). Rhesus monkeys know when they remember. *Proceedings of the National Academy of Sciences*, 98(9), 5359–5362.
- Hare, T. A., Camerer, C. F., & Rangel, A. (2009). Self-control in decision-making involves modulation of the vmPFC valuation system. *Science*, 324(5927), 646–648.
- Hart, J. T. (1965). Memory and the feeling-of-knowing experience. *Journal of educational psychology*, 56(4), 208.
- Heekeren, H. R., Marrett, S., Bandettini, P. A., & Ungerleider, L. G. (2004). A general mechanism for perceptual decision-making in the human brain. *Nature*, 431(7010), 859–862.
- Hertzog, C., Dixon, R. A., & Hultsch, D. F. (1990). Relationships between metamemory, memory predictions, and memory task performance in adults. *Psychology and aging*, 5(2), 215.
- Hoffrage, U. (2004). Overconfidence. In R. F. Pohl (Ed.), *Cognitive Illusions: A Handbook on Fallacies and Biases in Thinking, Judgement and Memory*, (pp. 235–254). Hove, UK: Psychology Press.
- Houser, D., & McCabe, K. (2014). Experimental Economics and Experimental Game Theory. In P. W. Glimcher & E. Fehr (Eds.), *Neuroeconomics* (2nd ed., pp. 19–34). Elsevier.
- Insabato, A., Pannunzi, M., Rolls, E. T., & Deco, G. (2010). Confidence-related decision making. *Journal of Neurophysiology*, 104(1), 539–547.
- Kaanders, P., Folke, T., & De Martino, B. (n.d.). Confidence informs information sampling and future decisions.
- Kahneman, D., & Tversky, A. (1979). Prospect theory: An analysis of decision under risk. *Econometrica: Journal of the Econometric Society*, 263–291.
- Kepecs, A., & Mainen, Z. F. (2012). A computational framework for the study of confidence in humans and animals. *Philosophical Transactions of the Royal Society of London B: Biological Sciences*, 367(1594), 1322–1337.
- Kepecs, A., & Mainen, Z. F. (2014). A Computational Framework for the Study of Confidence Across Species. In S. M. Fleming & C. D. Frith (Eds.), *The Cognitive Neuroscience of Metacognition*, (pp. 115–145). Springer.
- Kepecs, A., Uchida, N., Zariwala, H. A., & Mainen, Z. F. (2008). Neural correlates, computation and behavioural impact of decision confidence. *Nature*, 455(7210), 227–231.

- Keramati, M., Dezfouli, A., & Piray, P. (2011). Speed/accuracy trade-off between the habitual and the goal-directed processes. *PLoS Comput Biol*, 7(5), e1002055.
- Kiani, R., Corthell, L., & Shadlen, M. N. (2014). Choice Certainty Is Informed by Both Evidence and Decision Time. *Neuron*, 84(6), 1329–1342.
- Kiani, R., & Shadlen, M. N. (2009). Representation of confidence associated with a decision by neurons in the parietal cortex. *Science*, 324(5928), 759–764.
- Knill, D. C., & Pouget, A. (2004). The Bayesian brain: the role of uncertainty in neural coding and computation. *Trends in Neurosciences*, 27(12), 712–719.
- Kolling, N., Behrens, T. E. J., Mars, R. B., & Rushworth, M. F. S. (2012). Neural mechanisms of foraging. *Science*, 336(6077), 95–98.
- Komura, Y., Nikkuni, A., Hirashima, N., Uetake, T., & Miyamoto, A. (2013). Responses of pulvinar neurons reflect a subject's confidence in visual categorization. *Nature Neuroscience*, 16(6), 749–755.
- Koriat, A. (1993). How do we know that we know? The accessibility model of the feeling of knowing. *Psychological Review*, 100, 609–639.
- Koriat, A. (1997). Monitoring one's own knowledge during study: A cue-utilization approach to judgments of learning. *Journal of experimental psychology: general*, 126(4), 349.
- Koriat, A. (2007). Metacognition and consciousness. In P. Zelazo, M. Moscovitch, & E. Thompson (Eds.), *The Cambridge handbook of consciousness* (pp. 289–326). Cambridge University Press.
- Koriat, A., & Goldsmith, M. (1994). Memory in naturalistic and laboratory contexts: Distinguishing the accuracy-oriented and quantity-oriented approaches to memory assessment. *Journal of Experimental Psychology: General*, 123(3), 297.
- Koriat, A., & Goldsmith, M. (1996). Monitoring and control processes in the strategic regulation of memory accuracy. *Psychological review*, 103(3), 490.
- Koriat, A., Goldsmith, M., & Pansky, A. (2000). Toward a psychology of memory accuracy. *Annual review of psychology*, 51(1), 481–537.
- Koriat, A., & Levy-Sadot, R. (2001). The combined contributions of the cue-familiarity and accessibility heuristics to feelings of knowing. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 27(1), 34.
- Kornell, N., Son, L. K., & Terrace, H. S. (2007). Transfer of metacognitive skills and hint seeking in monkeys. *Psychological Science*, 18(1), 64–71.
- Krajbich, I., Armel, C., & Rangel, A. (2010). Visual fixations and the computation and comparison of value in simple choice. *Nature Neuroscience*, 13(10), 1292–1298.
- Krajbich, I., & Rangel, A. (2011). Multialternative drift-diffusion model predicts the relationship between visual fixations and choice in value-based decisions. *Proceedings of the National Academy of Sciences*, 108(33), 13852–13857.
- Kunimoto, C., Miller, J., & Pashler, H. (2001). Confidence and accuracy of near-threshold discrimination responses. *Consciousness and Cognition*, 10(3), 294–340.
- Lak, A., Costa, G. M., Romberg, E., Koulakov, A. A., Mainen, Z. F., & Kepecs, A. (2014). Orbitofrontal Cortex Is Required for Optimal Waiting Based on Decision Confidence.

- Neuron*, 84(1), 190–201.
- Lau, H. C., & Passingham, R. E. (2006). Relative blindsight in normal observers and the neural correlate of visual consciousness. *Proceedings of the National Academy of Sciences*, 103(49), 18763–18768.
- Lau, H., & Rosenthal, D. (2011). Empirical support for higher-order theories of conscious awareness. *Trends in Cognitive Sciences*, 15(8), 365–373.
- Lebreton, M., Abitbol, R., Daunizeau, J., & Pessiglione, M. (2015). Automatic integration of confidence in the brain valuation signal. *Nature Neuroscience*, 18(8), 1159–1167.
- Lebreton, M., Jorge, S., Michel, V., Thirion, B., & Pessiglione, M. (2009). An automatic valuation system in the human brain: evidence from functional neuroimaging. *Neuron*, 64(3), 431–439.
- Lee, M. D. (2008). BayesSDT: Software for Bayesian inference with signal detection theory. *Behavior Research Methods*, 40(2), 450–456.
- Lee, M. D., & Wagenmakers, E.-J. (2014). *Bayesian cognitive modeling: A practical course*. Cambridge university press.
- Levy, D. J., & Glimcher, P. W. (2012). The root of all value: a neural common currency for choice. *Current Opinion in Neurobiology*, 22(6), 1027–1038.
- Levy, I., Lazzaro, S. C., Rutledge, R. B., & Glimcher, P. W. (2011). Choice from non-choice: predicting consumer preferences from blood oxygenation level-dependent signals obtained during passive viewing. *Journal of Neuroscience*, 31(1), 118–125.
- Lim, S.-L., O'Doherty, J. P., & Rangel, A. (2011). The decision value computations in the vmPFC and striatum use a relative value code that is guided by visual attention. *Journal of Neuroscience*, 31(37), 13214–13223.
- Loomes, G., Starmer, C., & Sugden, R. (1991). Observing Violations of Transitivity by Experimental Methods. *Econometrica*, 59(2), 425–439.
- Louie, K., Khaw, M. W., & Glimcher, P. W. (2013). Normalization is a general neural mechanism for context-dependent decision making. *Proceedings of the National Academy of Sciences*, 110(15), 6139–6144.
- Ma, W. J., Beck, J. M., Latham, P. E., & Pouget, A. (2006). Bayesian inference with probabilistic population codes. *Nature Neuroscience*, 9(11), 1432–1438.
- MacIntyre, P. D., Noels, K. A., & Clément, R. (1997). Biases in self-ratings of second language proficiency: The role of language anxiety. *Language Learning*, 47(2), 265–287.
- Maniscalco, B., & Lau, H. (2012). A signal detection theoretic approach for estimating metacognitive sensitivity from confidence ratings. *Consciousness and Cognition*, 21(1), 422–430.
- Maniscalco, B., & Lau, H. (2014). Signal Detection Theory Analysis of Type 1 and Type 2 Data: Meta-d', Response-Specific Meta-d', and the Unequal Variance SDT Model (pp. 25–66). Springer.
- Maniscalco, B., & Lau, H. (2016). The signal processing architecture underlying subjective reports of sensory awareness. *Neuroscience of Consciousness*, 2016(1), niw002.
- Maniscalco, B., Peters, M. A. K., & Lau, H. (2016). Heuristic use of perceptual evidence leads to dissociation between performance and metacognitive sensitivity. *Attention, Perception, & Psychophysics*, 78(3), 923–937.

- Martin-Rhee, M. M., & Bialystok, E. (2008). The development of two types of inhibitory control in monolingual and bilingual children. *Bilingualism: Language and Cognition*, 11(1), 81–93.
- Masson, M. E. J., & Rotello, C. M. (2009). Sources of bias in the Goodman–Kruskal gamma coefficient measure of association: Implications for studies of metacognitive processes. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 35(2), 509–527.
- McCurdy, L. Y., Maniscalco, B., Metcalfe, J., Liu, K. Y., de Lange, F. P., & Lau, H. (2013). Anatomical coupling between distinct metacognitive systems for memory and visual perception. *Journal of Neuroscience*, 33(5), 1897–1906.
- Metcalfe, J., Van Snellenberg, J. X., DeRosse, P., Balsam, P., & Malhotra, A. K. (2014). Judgments of agency in schizophrenia: an impairment in autonoetic metacognition. In S. M. Fleming & C. D. Frith (Eds.) *The Cognitive Neuroscience of Metacognition* (pp. 367–387). Springer.
- Meyniel, F., & Dehaene, S. (2017). Brain networks for confidence weighting and hierarchical inference during probabilistic learning. *Proceedings of the National Academy of Sciences*, 114(19), E3859–E3868.
- Meyniel, F., Schlunegger, D., & Dehaene, S. (2015). The sense of confidence during probabilistic learning: A normative account. *PLoS Comput Biol*, 11(6), e1004305.
- Meyniel, F., Sigman, M., & Mainen, Z. F. (2015). Confidence as Bayesian Probability: From Neural Origins to Behavior. *Neuron*, 88(1), 78–92.
- Moran, R., Teodorescu, A. R., & Usher, M. (2015). Post choice information integration as a causal determinant of confidence: Novel data and a computational account. *Cognitive Psychology*, 78, 99–147.
- Morey, R. D., & Rouder, J. N. (2011). Bayes factor approaches for testing interval null hypotheses. *Psychological Methods*, 16(4), 406–419.
- Morton, J. B., & Harper, S. N. (2007). What did Simon say? Revisiting the bilingual advantage. *Developmental Science*, 10(6), 719–26.
- Navajas, J., Bahrami, B., & Latham, P. E. (2016). Post-decisional accounts of biases in confidence. *Current Opinion in Behavioral Sciences*, 11, 55–60.
- Navajas, J., Hindocha, C., Foda, H., Keramati, M., Latham, P. E., & Bahrami, B. (2017). The idiosyncratic nature of confidence. *bioRxiv*, 102269. Now in Nature Hum Behaviour
- Nelson, T. O. (1984). A comparison of current measures of the accuracy of feeling-of-knowing predictions. *Psychological Bulletin*, 95(1), 109.
- Nelson, T. O. (1996). Consciousness and metacognition. *American Psychologist*, 51(2), 102.
- Nelson, T. O., & Narens, L. (1990). Metamemory: A theoretical framework and new findings. In G. Bower (Ed.), *The psychology of learning and motivation: Advances in research and theory* (Vol. 26, pp. 125–173). New York: Academic Press.
- Nelson, T. O., & Narens, L. (1994). Why investigate metacognition? In J. Metcalfe & A. P. Shimamura (Eds.), *Metacognition: Knowing about knowing* (pp. 1–25). Cambridge, MA: MIT Press.
- Nieuwenhuis, S., Ridderinkhof, K. R., Blom, J., Band, G. P. H., & Kok, A. (2001). Error-related brain potentials are differentially related to awareness of response errors: Evidence from an

- antisaccade task. *Psychophysiology*, 38(5), 752–760.
- Notebaert, W., Houtman, F., Van Opstal, F., Gevers, W., Fias, W., & Verguts, T. (2009). Post-error slowing: an orienting account. *Cognition*, 111(2), 275–279.
- Paap, K. R., Johnson, H. A., & Sawi, O. (2015). Bilingual advantages in executive functioning either do not exist or are restricted to very specific and undetermined circumstances. *Cortex*, 69, 265–278.
- Payzan-LeNestour, E., Dunne, S., Bossaerts, P., & O’Doherty, J. P. (2013). The neural representation of unexpected uncertainty during value-based decision making. *Neuron*, 79(1), 191–201.
- Pedings, K. E., Langville, A. N., & Yamamoto, Y. (2012). A minimum violations ranking method. *Optimization and Engineering*, 13(2), 349–370.
- Peirce, C. S., & Jastrow, J. (1884). On small differences of sensation. *Mem. Natl. Acad. Sci.*, (3), 75–83.
- Peirce, J. W. (2008). Generating stimuli for neuroscience using PsychoPy. *Frontiers in neuroinformatics*, 2, 10.
- Plassmann, H., O’Doherty, J. P., & Rangel, A. (2010). Appetitive and aversive goal values are encoded in the medial orbitofrontal cortex at the time of decision making. *Journal of Neuroscience*, 30(32), 10799–10808.
- Pleskac, T. J., & Busemeyer, J. R. (2010). Two-stage dynamic signal detection: a theory of choice, decision time, and confidence. *Psychological Review*, 117(3), 864–901.
- Pouget, A., Drugowitsch, J., & Kepecs, A. (2016). Confidence and certainty: distinct probabilistic quantities for different goals. *Nat Neurosci*, 19(3), 366–374.
- Rabbitt, P., Cumming, G., & Vyas, S. (1978). Some errors of perceptual analysis in visual search can be detected and corrected. *The Quarterly Journal of Experimental Psychology*, 30(2), 319–332.
- Rabbitt, P. M. (1966). Errors and error correction in choice-response tasks. *Journal of Experimental Psychology*, 71(2), 264–272.
- Ratcliff, R. (1978). A theory of memory retrieval. *Psychological Review*, 85(2), 59–108.
- Ratcliff, R. (2002). A diffusion model account of response time and accuracy in a brightness discrimination task: Fitting real data and failing to fit fake but plausible data. *Psychonomic Bulletin & Review*, 9(2), 278–291.
- Ratcliff, R., & Childers, R. (2015). Individual differences and fitting methods for the two-choice diffusion model of decision making. *Decision*, 2(4), 237–279.
- Ratcliff, R., & McKoon, G. (2008). The diffusion decision model: theory and data for two-choice decision tasks. *Neural Computation*, 20(4), 873–922.
- Ratcliff, R., & Smith, P. L. (2004). A comparison of sequential sampling models for two-choice reaction time. *Psychological Review*, 111(2), 333–367.
- Ratcliff, R., & Starns, J. J. (2009). Modeling confidence and response time in recognition memory. *Psychological Review*, 116(1), 59–83.
- Ratcliff, R., Thapar, A., & McKoon, G. (2010). Individual differences, aging, and IQ in two-choice tasks. *Cognitive Psychology*, 60(3), 127–157.

- Ratcliff, R., & Tuerlinckx, F. (2002). Estimating parameters of the diffusion model: Approaches to dealing with contaminant reaction times and parameter variability. *Psychonomic bulletin & review*, 9(3), 438-481.
- Raven, J. C., & Court, J. H. (1986). Raven's progressive matrices and Raven's coloured matrices. *London: HK Lewis*.
- Reder, L. M. (1987). Strategy selection in question answering. *Cognitive psychology*, 19(1), 90-138.
- Resulaj, A., Kiani, R., Wolpert, D. M., & Shadlen, M. N. (2009). Changes of mind in decision-making. *Nature*, 461(7261), 263–266.
- Rhodes, M. G., & Tauber, S. K. (2011). The influence of delaying judgments of learning on metacognitive accuracy: A meta-analytic review. *Psychological bulletin*, 137(1), 131.
- Rounis, E., Maniscalco, B., Rothwell, J. C., Passingham, R. E., & Lau, H. (2010). Theta-burst transcranial magnetic stimulation to the prefrontal cortex impairs metacognitive visual awareness. *Cognitive Neuroscience*, 1(3), 165–175.
- Roxin, A., & Ledberg, A. (2008). Neurobiological models of two-choice decision making can be reduced to a one-dimensional nonlinear diffusion equation. *PLoS Comput Biol*, 4(3), e1000046.
- Rubio-Fernández, P., & Glucksberg, S. (2012). Reasoning about other people's beliefs: bilinguals have an advantage. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 38(1), 211-217.
- Rushworth, M. F., Kolling, N., Sallet, J., & Mars, R. B. (2012). Valuation and decision-making in frontal cortex: one or many serial or parallel systems? *Current Opinion in Neurobiology*, 22(6), 946–955.
- Sanders, J. I., Hangya, B., & Kepecs, A. (2016). Signatures of a statistical computation in the human sense of confidence. *Neuron*, 90(3), 499–506.
- Scheffers, M. K., & Coles, M. G. H. (2000). Performance monitoring in a confusing world: error-related brain activity, judgments of response accuracy, and types of errors. *Journal of Experimental Psychology: Human Perception and Performance*, 26(1), 141-151.
- Schwartz, B. L., & Díaz, F. (2014). Quantifying human metacognition for the neurosciences. In S. M. Fleming & C. D. Frith (Eds.), *The Cognitive Neuroscience of Metacognition*, (pp. 9–23). Springer Berlin Heidelberg.
- Shea, N., Boldt, A., Bang, D., Yeung, N., Heyes, C., & Frith, C. D. (2014). Supra-personal cognitive control and metacognition. *Trends in Cognitive Sciences*, 18(4), 186–193.
- Shimamura, A. P. (2000). Toward a cognitive neuroscience of meta-cognition. *Consciousness and Cognition*, 9, 313 – 323.
- Siedlecka, M., Paulewicz, B., & Wierzchoń, M. (2016). But I Was So Sure! Metacognitive Judgments Are Less Accurate Given Prospectively than Retrospectively. *Frontiers in Psychology*, 7, 218.
- Smith, J. D., Couchman, J. J., & Beran, M. J. (2014). The highs and lows of theoretical interpretation in animal-metacognition research. In *The Cognitive Neuroscience of Metacognition* (pp. 87-111). Springer, Berlin, Heidelberg.
- Smith J. D., Schull J., Strote J., McGee K., Egnor R., Erb L. (1995) The uncertain response in the bottlenosed dolphin (*Tursiops truncatus*). *J Exp Psychol Gen* 124:391–408.

- Soltani, A., De Martino, B., & Camerer, C. (2012). A range-normalization model of context-dependent choice: a new model and evidence. *PLoS Comput Biol*, 8(7), e1002607.
- Song, C., Kanai, R., Fleming, S. M., Weil, R. S., Schwarzkopf, D. S., & Rees, G. (2011). Relating inter-individual differences in metacognitive performance on different perceptual tasks. *Consciousness and Cognition*, 20(4), 1787–1792.
- Stanislaw, H., & Todorov, N. (1999). Calculation of signal detection theory measures. *Behavior Research Methods, Instruments, & Computers*, 31(1), 137–149.
- Stocco, A., & Prat, C. S. (2014). Bilingualism trains specific brain circuits involved in flexible rule selection and application. *Brain and Language*, 137, 50–61.
- Stocco, A., Yamasaki, B., Natalenko, R., & Prat, C. S. (2014). Bilingual brain training: A neurobiological framework of how bilingual experience improves executive function. *International Journal of Bilingualism*, 18(1), 67–92.
- Summerfield, C., & Tsetsos, K. (2012). Building bridges between perceptual and economic decision-making: neural and computational mechanisms. *Frontiers in Neuroscience*, 6, 70.
- Swets, J. A. (2014). *Signal detection theory and ROC analysis in psychology and diagnostics: Collected papers*. Mahwah, NJ: Lawrence Erlbaum Associates, Inc.
- Thompson, V. A., Turner, J. A. P., & Pennycook, G. (2011). Intuition, reason, and metacognition. *Cognitive psychology*, 63(3), 107–140.
- Thompson, V. A., Turner, J. A. P., Pennycook, G., Ball, L. J., Brack, H., Ophir, Y., & Ackerman, R. (2013). The role of answer fluency and perceptual fluency as metacognitive cues for initiating analytic thinking. *Cognition*, 128(2), 237–251.
- Treisman, M., & Faulkner, A. (1984). The setting and maintenance of criteria representing levels of confidence. *Journal of Experimental Psychology: Human Perception and Performance*, 10(1), 119–139.
- Trofimovich, P., Isaacs, T., Kennedy, S., Saito, K., & Crowther, D. (2016). Flawed self-assessment: Investigating self-and other-perception of second language speech. *Bilingualism: Language and Cognition*, 19(1), 122–140.
- Trommershäuser, J., Maloney, L. T., & Landy, M. S. (2003). Statistical decision theory and trade-offs in the control of motor response. *Spatial Vision*, 16(3), 255–275.
- Van Den Berg, R., Anandalingam, K., Zylberberg, A., Kiani, R., Shadlen, M. N., & Wolpert, D. M. (2016). A common mechanism underlies changes of mind about decisions and confidence. *eLife*, 5, e12192.
- Van den Berg, R., Zylberberg, A., Kiani, R., Shadlen, M. N., & Wolpert, D. M. (2016). Confidence Is the Bridge between Multi-stage Decisions. *Current Biology*, 26(23), 3157–3168.
- Van Heuven, W. J. B., Schriefers, H., Dijkstra, T., & Hagoort, P. (2008). Language conflict in the bilingual brain. *Cerebral Cortex*, 18(11), 2706–2716.
- Vandekerckhove, J., Tuerlinckx, F., & Lee, M. D. (2011). Hierarchical diffusion models for two-choice response times. *Psychological Methods*, 16(1), 44–62.
- Veenman, M. V. J., Elshout, J. J., & Meijer, J. (1997). The generality vs domain-specificity of metacognitive skills in novice learning across domains. *Learning and Instruction*, 7(2), 187–

- Vickers, D., & Packer, J. (1982). Effects of alternating set for speed or accuracy on response time, accuracy and confidence in a unidimensional discrimination task. *Acta Psychologica*, 50(2), 179–197.
- Vlassova, A., Donkin, C., & Pearson, J. (2014). Unconscious information changes decision accuracy but not confidence. *Proceedings of the National Academy of Sciences*, 111(45), 16214–16218.
- Von Neumann, J., & Morgenstern, O. (2007). *Theory of games and economic behavior*. Princeton university press.
- Voss, A., Rothermund, K., & Voss, J. (2004). Interpreting the parameters of the diffusion model: An empirical validation. *Memory & Cognition*, 32(7), 1206–1220.
- Washburn, D. A., Smith, J. D., & Shields, W. E. (2006). Rhesus monkeys (*Macaca mulatta*) immediately generalize the uncertain response. *Journal of Experimental Psychology: Animal Behavior Processes*, 32(2), 185.
- Wechsler, D. (2008). (WAIS-IV). *Administration and Scoring Manual*. San Antonio, TX, USA: The Psychological Corporation.
- Wei, Z., & Wang, X.-J. (2015). Confidence estimation as a stochastic process in a neurodynamical system of decision making. *Journal of Neurophysiology*, 114(1), 99–113.
- Yeung, N., & Summerfield, C. (2012). Metacognition in human decision-making: confidence and error monitoring. *Phil. Trans. R. Soc. B*, 367(1594), 1310–1321.
- Yeung, N., & Summerfield, C. (2014). Shared Mechanisms for Confidence Judgements and Error Detection in Human Decision Making. In S. M. Fleming & C. D. Frith (Eds.), *The Cognitive Neuroscience of Metacognition*, (pp. 147–167). Springer.
- Yoshida, W., & Ishii, S. (2006). Resolution of uncertainty in prefrontal cortex. *Neuron*, 50(5), 781–789.
- Yu, S., Pleskac, T. J., & Zeigenfuse, M. D. (2015). Dynamics of postdecisional processing of confidence. *Journal of Experimental Psychology: General*, 144(2), 489.
- Zylberberg, A., Barttfeld, P., & Sigman, M. (2012). The construction of confidence in a perceptual decision. *Frontiers in Integrative Neuroscience*, 6, 79.

Appendices

Appendix 1: List of stimuli in Experiment 3:

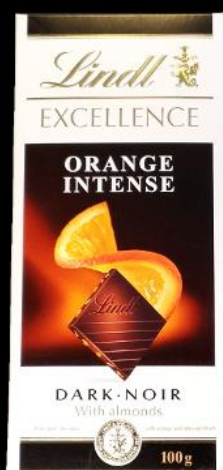






Appendix 2: List of stimuli in Experiment 4:

























Appendix 3: GSF Does Not Predict Choice, but Interacts With Stimulus Strength in a Model that does not Include DDT

A hierarchical logistic regression model with participantwise variation in intercepts and slopes predicted the choice of the reference item (the first item encountered in western reading order) from the difference in value between the reference item and the mean value of the other two items, GSF and a value difference GSF interaction. Value difference strongly predicted choice ($z=8.59$, $p<10^{-10}$), as did the interaction term ($z=-4.26$, $p<10^{-4}$) but GSF alone was not a significant predictor of choice ($z=1.42$, $p=.16$).

Appendix 4: Confidence Distributions for Simultaneous and Sequential Confidence Judgments in Experiment 6

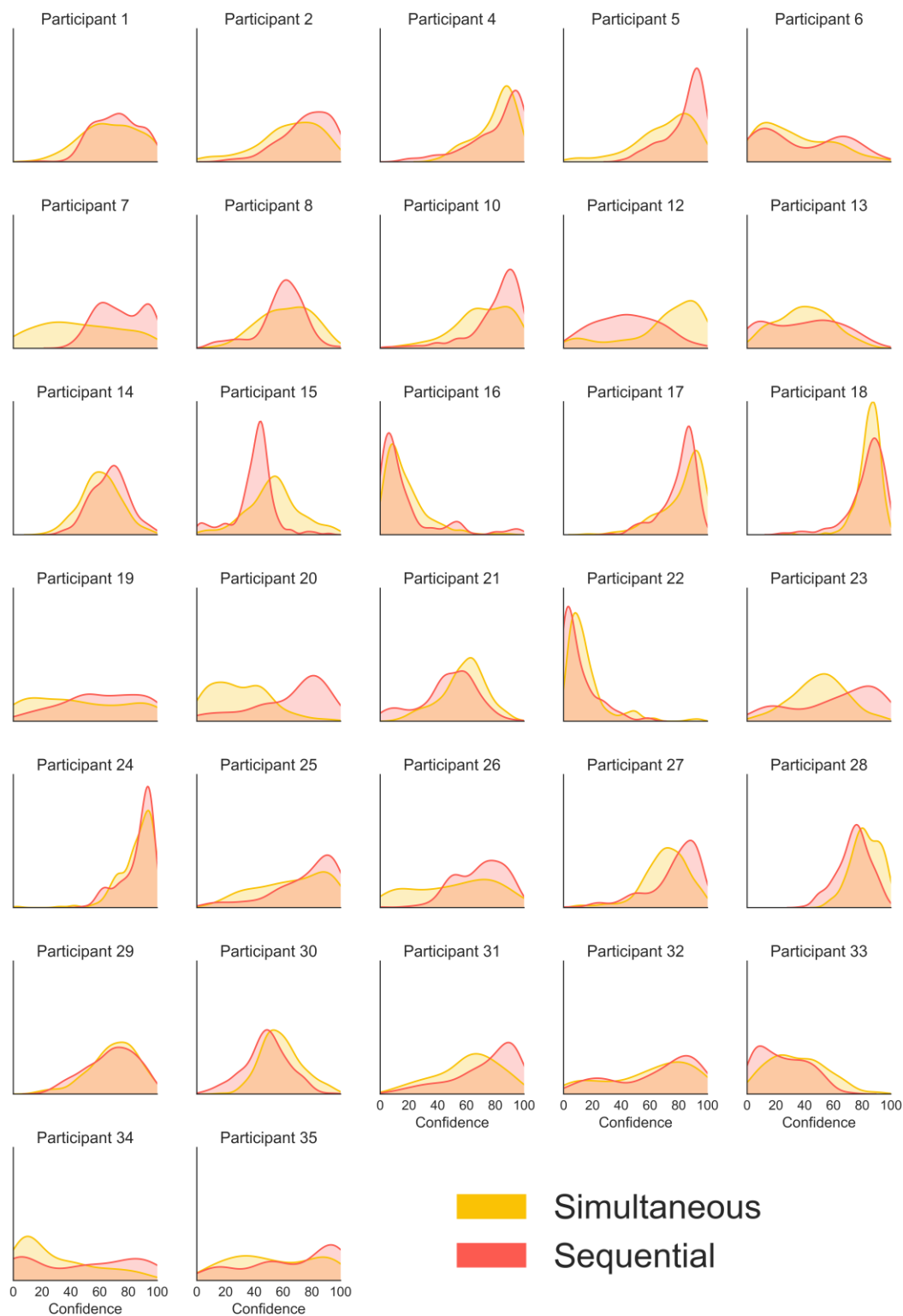


Figure A1. The Response Distributions for Sequentially and Simultaneously Reported Confidence

Appendix 5: Interjudgment Times Are Not Associated with Confidence, Accuracy or Stimulus Strength in Experiment 6

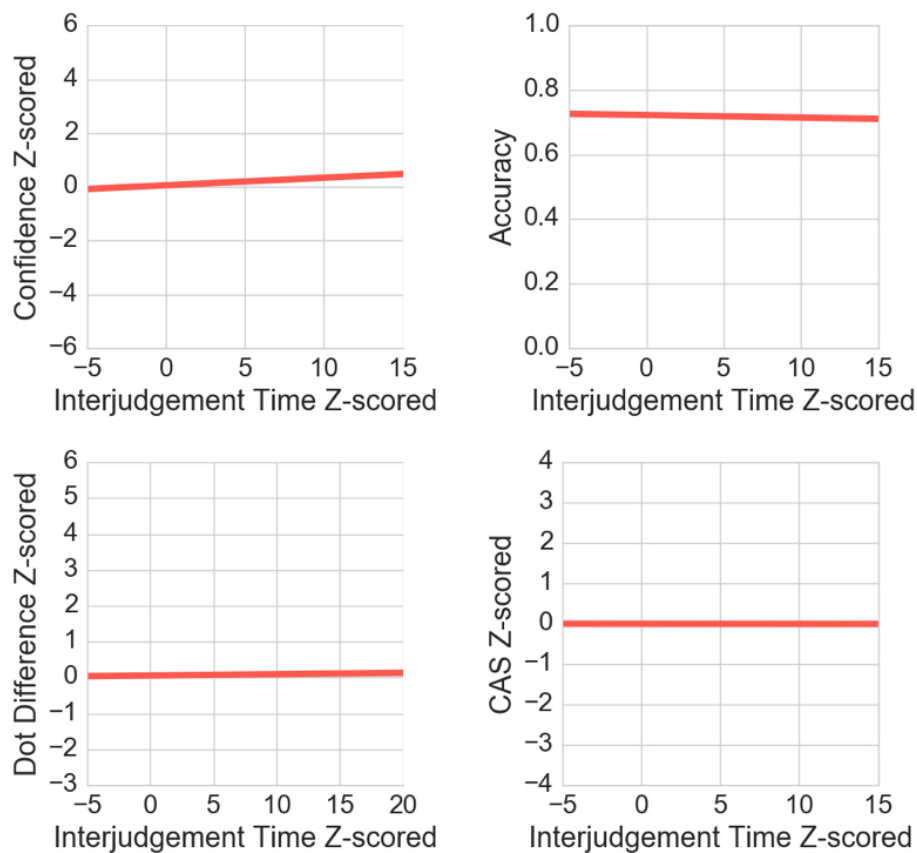


Figure A2. Relationship Between Interjudgment Times and Confidence, First-order Accuracy, Stimulus strength and Second-order Accuracy.

I ran a set of hierarchical regression models with varying slopes and intercepts to test whether interjudgment time interacted with any of the other variables. All variables reported in these analyses were z-scored. Interjudgment times did not predict the magnitude of reported confidence ($t=1.31$, $p=0.16$), nor were they associated with the accuracy of the choices ($z=-0.09$, $p=.93$), nor were they associated with the stimulus strength of the trial (here dot difference; $t=0.14$, $p=.88$). To test whether interjudgment times were related to second order accuracy I created a confidence accuracy score (CAS): Confidence * Correct. Where confidence ranged between one and 100 and correct was set to 1 for correct trials and -1 for incorrect trials. Interjudgment time did not predict CAS ($t=-0.01$, $p=.99$). For reference, RT did predict CAS ($t=-12.68$, $p<10^{-10}$), as did stimulus strength ($t=20.17$, $p<10^{-10}$).

Appendix 6: Confidence is Influenced by Both Positive and Negative Evidence for Both Sequential and Simultaneous Confidence Judgments in Experiment 6

To evaluate whether both positive and negative evidence influenced confidence I ran two hierarchical linear regression models with intercepts varying by participants. Both positive ($t=11.17, p=10^{-10}$) and negative evidence ($t=-8.99, p<10^{-10}$) predicted confidence for the simultaneous confidence responses. And both positive ($t=14.09, p=10^{-10}$) and negative evidence ($t=-13.95, p<10^{-10}$) predicted confidence for the sequential confidence responses. Therefore, while a model with separate weight parameters for positive and negative evidence fitted the data better than a model that only included the difference ($BIC_{\text{Difference Model}}=44\ 118$ $BIC_{\text{Separate Weights Model}}=43\ 747$), I opted for the simpler representation of stimulus strength, given the complexity of some of the models in Chapter 5.